

UNITED STATES PATENT APPLICATION

TITLE: METHODS FOR MODIFYING PLANT BIOMASS AND TOLERANCE TO ABIOTIC STRESS

INVENTORS: **JIANG, Cai-Zhong**
 HEARD, Jacqueline E.
 RATCLIFFE, Oliver
 GUTTERSON, Neal
 HEMPEL, Frederick
 KUMIMOTO, Roderick W.
 KEDDIE, James S.
 SHERMAN, Bradley K.

CERTIFICATE OF EXPRESS MAILING

"Express Mail" Label No.: EL 889 523 626 US

Date of Deposit: September 23, 2003

I HEREBY CERTIFY UNDER 37 C.F.R. 1.10 THAT THIS CORRESPONDENCE IS BEING DEPOSITED WITH THE UNITED STATES POSTAL SERVICE AS "EXPRESS MAIL POST OFFICE TO ADDRESSEE" WITH SUFFICIENT POSTAGE ON THE DATE INDICATED ABOVE AND IS ADDRESSED TO :

MAIL STOP PATENT APPLICATION

COMMISSIONER FOR PATENTS

PO BOX 1450

ALEXANDRIA, VA 22313-1450

Kathleen K. Muto
(Signature)

KATHLEEN MUTO
(Printed Name)

METHODS FOR MODIFYING PLANT BIOMASS AND TOLERANCE TO ABIOTIC STRESS

RELATIONSHIP TO COPENDING APPLICATIONS

5

This application claims the benefit of copending US Non-provisional Application No. 10/374,780, filed February 25, 2003; which claims the benefit of US Provisional Application No. 60/336,049, filed November 19, 2001, US Non-provisional Application No. 09/934,455, filed August 22, 2001, which in turn claims priority from US Provisional Application No. 60/227,439, filed August 22, 2000, and US Provisional Application No. 60/310,847, filed August 9, 2001; US Non-provisional Application No. 10/412,699, filed April 10, 2003; which claims the benefit of US Non-provisional Application No. 09/506,720, filed February 17, 2000, which in turn claims the benefit of US Provisional Application No. 60/135,134, filed May 20, 1999, US Non-provisional Application No. 09/533,392, filed March 22, 2000, US Non-provisional Application No. 09/533,029, filed March 22, 2000, US Non-provisional Application No. 09/532,591, filed March 22, 2000, which in turn claimed the benefit of US Provisional Application No. 60/125,814, filed March 23, 1999, US Non-provisional Application No. 09/533,030, filed March 22, 2000, US Non-provisional Application No. 09/713,994, filed November 16, 2000, US Non-provisional Application No. 09/996,140, filed November 26, 2001, US Non-provisional Application No. 09/823,676, filed April 2, 2001; US Non-provisional Application No. 10/421,138, filed April 23, 2003; US Non-provisional Application No. 10/225,068, filed August 9, 2002, copending US Non-provisional Application No. 10/225,066, filed August 9, 2002, copending US Non-provisional Application No. 10/225,067, filed August 9, 2002, filed August 9, 2002, which claim the benefit of US Provisional Application No. 60/338,692, filed December 11, 2001. The entire contents of these applications are hereby incorporated by reference.

25

FIELD OF THE INVENTION

The present invention relates to polynucleotides comprising plant genes or fragments of plant genes that increase a plant's size or biomass, the yield that may be obtained from such a plant, and compositions and methods for producing plants having increased size or biomass. The invention also pertains to plants having altered sugar sensing and increased tolerance to abiotic stresses, including osmotic stresses such as drought, salt stress, heat stress, and germination in cold conditions.

30

BACKGROUND OF THE INVENTION

A plant's traits, such as its biochemical, developmental, or phenotypic characteristics, may be controlled through a number of cellular processes. One important way to manipulate that control is through transcription factors – proteins that influence the expression of a particular gene or sets of genes, for example, those that affect a plant's size or tolerance to abiotic stresses. Transformed and transgenic plants that comprise cells having altered levels of at least one selected transcription factor, for example, possess advantageous or desirable traits. Strategies for manipulating traits by altering a plant cell's transcription factor content can therefore result in plants and crops with new and/or improved commercially valuable properties.

Transcription factors can modulate gene expression, either increasing or decreasing (inducing or repressing) the rate of transcription. This modulation results in differential levels of gene expression at various developmental stages, in different tissues and cell types, and in response to different exogenous (e.g., environmental) and endogenous stimuli throughout the life cycle of the organism.

Phylogenetic relationships among organisms have been demonstrated many times, and studies from a diversity of prokaryotic and eukaryotic organisms suggest a more or less gradual evolution of biochemical and physiological mechanisms and metabolic pathways. Despite different evolutionary pressures, proteins that regulate the cell cycle in yeast, plant, nematode, fly, rat, and man have common chemical or structural features and modulate the same general cellular activity. Comparisons of *Arabidopsis* gene sequences with those from other organisms where the structure and/or function may be known allow researchers to draw analogies and to develop model systems for testing hypotheses. These model systems are of great importance in developing and testing plant varieties with novel traits that may have an impact upon agronomy.

Because transcription factors are key controlling elements of biological pathways, altering the expression levels of one or more transcription factors can change entire biological pathways in an organism. For example, manipulation of the levels of selected transcription factors may result in increased expression of economically useful proteins or biomolecules in plants or improvement in other agriculturally relevant characteristics. Conversely, blocked or reduced expression of a transcription factor may reduce biosynthesis of unwanted compounds or remove an undesirable trait. Therefore, manipulating transcription factor levels in a plant offers tremendous potential in agricultural biotechnology for modifying a plant's traits, including traits that improve yield, or a plant's survival and yield during periods of abiotic stress, including, for example, germination in cold conditions, excessive heat, and osmotic stresses such as drought and salt stress.

Desirability of increasing biomass. The ability to increase the biomass or size of a plant would have several important commercial applications. Crop species may be generated that produce higher yields on larger cultivars, particularly those in which the vegetative portion of the plant is edible. For example, increasing plant leaf biomass may increase the yield of leafy vegetables for human or animal consumption. Additionally, increasing leaf biomass can be used to increase production of plant-derived pharmaceutical or industrial products. By increasing plant biomass, increased production levels of the products may be obtained from the plants. Tobacco leaves, in particular, have been employed as plant factories to generate such products. Furthermore, it may be desirable to increase crop yields of plants by increasing total plant photosynthesis. An increase in total plant photosynthesis is typically achieved by increasing leaf area of the plant. Additional photosynthetic capacity may be used to increase the yield derived from particular plant tissue, including the leaves, roots, fruits or seed. In addition, the ability to modify the biomass of the leaves may be useful for permitting the growth of a plant under decreased light intensity or under high light intensity. Modification of the biomass of another tissue, such as roots, may be useful to improve a plant's ability to grow under harsh environmental conditions, including drought or nutrient deprivation, because the roots may grow deeper into the ground. Increased biomass can also be a consequence of some strategies for increased tolerance to stresses, such as drought stress. Early in a stress response plant growth (e.g., expansion of lateral organs, increase in stem girth, etc.) can be slowed to enable the plant to activate adaptive responses. Growth rate that is less sensitive to stress-induced control can result in enhanced plant size, particularly later in development.

For some ornamental plants, the ability to provide larger varieties would be highly desirable. For many plants, including fruit-bearing trees, trees that are used for lumber production, or trees and shrubs that serve as view or wind screens, increased stature provides improved benefits in the forms of greater yield or improved screening.

Because increased yield may be quite valuable to growers, we believe that there is significant commercial opportunity for engineering pathogen tolerance or resistance using transgenic plants with altered expression of the instant plant transcription factors. Crops so engineered will provide higher yields, and may be used to improve the appearance of ornamentals. The present invention satisfies a need in the art by providing new compositions that are useful for engineering plants with increased biomass or size, and having the potential to increase yield.

Problems associated with drought. A drought is a period of abnormally dry weather that persists long enough to produce a serious hydrologic imbalance (for example crop damage, water supply shortage, etc.). While much of the weather that we experience is brief and short-lived, drought is a more gradual phenomenon, slowly taking hold of an area and tightening its grip with time. In severe cases, drought can

last for many years and can have devastating effects on agriculture and water supplies. With burgeoning population and chronic shortage of available fresh water, drought is not only the number one weather related problem in agriculture, it also ranks as one of the major natural disasters of all time, causing not only economic damage, but also loss of human lives. For example, losses from the US drought of 1988 exceeded \$40 billion, exceeding the losses caused by Hurricane Andrew in 1992, the Mississippi River floods of 1993, and the San Francisco earthquake in 1989. In some areas of the world, the effects of drought can be far more severe. In the Horn of Africa the 1984–1985 drought led to a famine that killed 750,000 people.

Problems for plants caused by low water availability include mechanical stresses caused by the withdrawal of cellular water. Drought also causes plants to become more susceptible to various diseases (Simpson (1981). "The Value of Physiological Knowledge of Water Stress in Plants", In Water Stress on Plants, (Simpson, G. M., ed.), Praeger, NY, pp. 235-265).

In addition to the many land regions of the world that are too arid for most if not all crop plants, overuse and over-utilization of available water is resulting in an increasing loss of agriculturally-usable land, a process which, in the extreme, results in desertification. The problem is further compounded by increasing salt accumulation in soils, as described above, which adds to the loss of available water in soils.

Problems associated with high salt levels. One in five hectares of irrigated land is damaged by salt, an important historical factor in the decline of ancient agrarian societies. This condition is only expected to worsen, further reducing the availability of arable land and crop production, since none of the top five food crops - wheat, corn, rice, potatoes, and soybean - can tolerate excessive salt.

Detrimental effects of salt on plants are a consequence of both water deficit resulting in osmotic stress (similar to drought stress) and the effects of excess sodium ions on critical biochemical processes. As with freezing and drought, high saline causes water deficit; the presence of high salt makes it difficult for plant roots to extract water from their environment (Buchanan et al. (2000) in Biochemistry and Molecular Biology of Plants, American Society of Plant Physiologists, Rockville, MD). Soil salinity is thus one of the more important variables that determines where a plant may thrive. In many parts of the world, sizable land areas are uncultivable due to naturally high soil salinity. To compound the problem, salination of soils that are used for agricultural production is a significant and increasing problem in regions that rely heavily on agriculture. The latter is compounded by over-utilization, over-fertilization and water shortage, typically caused by climatic change and the demands of increasing population. Salt tolerance is of particular importance early in a plant's lifecycle, since evaporation from the soil surface causes upward water movement, and salt accumulates in the upper soil layer where the seeds are placed. Thus, germination normally takes place at a salt concentration much higher than the mean salt level in the

whole soil profile.

Problems associated with excessive heat. Germination of many crops is very sensitive to temperature. A transcription factor that would enhance germination in hot conditions would be useful for crops that are planted late in the season or in hot climates. Seedlings and mature plants that are exposed to excess heat may experience heat shock, which may arise in various organs, including leaves and particularly fruit, when transpiration is insufficient to overcome heat stress. Heat also damages cellular structures, including organelles and cytoskeleton, and impairs membrane function (Buchanan et al. (2000) in Biochemistry and Molecular Biology of Plants, American Society of Plant Physiologists, Rockville, MD).

Heat shock may produce a decrease in overall protein synthesis, accompanied by expression of heat shock proteins. Heat shock proteins function as chaperones and are involved in refolding proteins denatured by heat.

Heat stress often accompanies conditions of low water availability. Heat itself is seen as an interacting stress and adds to the detrimental effects caused by water deficit conditions. Evaporative demand exhibits near exponential increases with increases in daytime temperatures and can result in high transpiration rates and low plant water potentials (Hall et al. (2000) *Plant Physiol.* 123: 1449-1458). High-temperature damage to pollen almost always occurs in conjunction with drought stress, and rarely occurs under well-watered conditions. Thus, separating the effects of heat and drought stress on pollination is difficult. Combined stress can alter plant metabolism in novel ways; therefore understanding the interaction between different stresses may be important for the development of strategies to enhance stress tolerance by genetic manipulation.

Problems associated with excessive chilling conditions. The term "chilling sensitivity" has been used to describe many types of physiological damage produced at low, but above freezing, temperatures. Most crops of tropical origins such as soybean, rice, maize and cotton are easily damaged by chilling. Typical chilling damage includes wilting, necrosis, chlorosis or leakage of ions from cell membranes. The underlying mechanisms of chilling sensitivity are not completely understood yet, but probably involve the level of membrane saturation and other physiological deficiencies. For example, photoinhibition of photosynthesis (disruption of photosynthesis due to high light intensities) often occurs under clear atmospheric conditions subsequent to cold late summer/autumn nights. For example, chilling may lead to yield losses and lower product quality through the delayed ripening of maize. Another consequence of poor growth is the rather poor ground cover of maize fields in spring, often resulting in soil erosion, increased occurrence of weeds, and reduced uptake of nutrients. A retarded uptake of mineral nitrogen could also lead to increased losses of nitrate into the ground water. By some estimates, chilling accounts

for monetary losses in the United States (US) behind only to drought and flooding.

Desirability of altered sugar sensing. Sugars are key regulatory molecules that affect diverse processes in higher plants including germination, growth, flowering, senescence, sugar metabolism and photosynthesis. Sucrose, for example, is the major transport form of photosynthate and its flux through cells has been shown to affect gene expression and alter storage compound accumulation in seeds (source-sink relationships). Glucose-specific hexose-sensing has also been described in plants and is implicated in cell division and repression of "famine" genes (photosynthetic or glyoxylate cycles).

Water deficit is a common component of many plant stresses. Water deficit occurs in plant cells when the whole plant transpiration rate exceeds the water uptake. In addition to drought, other stresses, such as salinity and low temperature, produce cellular dehydration (McCue and Hanson (1990) *Trends Biotechnol.* 8: 358-362).

Salt and drought stress signal transduction consist of ionic and osmotic homeostasis signaling pathways. The ionic aspect of salt stress is signaled via the SOS pathway where a calcium-responsive SOS3-SOS2 protein kinase complex controls the expression and activity of ion transporters such as SOS1. The pathway regulating ion homeostasis in response to salt stress has been reviewed recently by Xiong and Zhu (Xiong and Zhu (2002) *Plant Cell Environ.* 25: 131-139).

The osmotic component of salt stress involves complex plant reactions that overlap with drought and/or cold stress responses.

Common aspects of drought, cold and salt stress response have been reviewed recently by Xiong and Zhu (2002) *supra*. Those include:

- (a) transient changes in the cytoplasmic calcium levels very early in the signaling event (Knight, (2000) *Int. Rev. Cytol.* 195: 269-324; Sanders et al. (1999) *Plant Cell* 11: 691-706);
- (b) signal transduction via mitogen-activated and/or calcium dependent protein kinases (CDPKs; see Xiong and Zhu (2002) *supra*) and protein phosphatases (Merlot et al. (2001) *Plant J.* 25: 295-303; Tähtiharju and Palva (2001) *Plant J.* 26: 461-470);
- (c) increases in abscisic acid levels in response to stress triggering a subset of responses (Xiong and Zhu (2002) *supra*, and references therein);
- (d) inositol phosphates as signal molecules (at least for a subset of the stress responsive transcriptional changes (Xiong et al. (2001) *Genes Dev.* 15: 1971-1984);
- (e) activation of phospholipases which in turn generate a diverse array of second messenger molecules, some of which might regulate the activity of stress responsive kinases (phospholipase D functions in an ABA independent pathway, Frank et al. (2000) *Plant Cell* 12: 111-124);
- (f) induction of late embryogenesis abundant (LEA) type genes including the CRT/DRE-

containing COR/RD genes (Xiong and Zhu (2002) *supra*);

- (g) increased levels of antioxidants and compatible osmolytes such as proline and soluble sugars (Hasegawa et al. (2000) *Annu. Rev. Plant Mol. Plant Physiol.* 51: 463-499);
- (h) accumulation of reactive oxygen species such as superoxide, hydrogen peroxide, and hydroxyl radicals (Hasegawa et al. (2000) *supra*).

Abscisic acid biosynthesis is regulated by osmotic stress at multiple steps. Both ABA-dependent and ABA-independent osmotic stress signaling first modify constitutively expressed transcription factors, leading to the expression of early response transcriptional activators, which then activate downstream stress tolerance effector genes.

Based on the commonality of many aspects of cold, drought and salt stress responses, it can be concluded that genes that increase tolerance to cold or salt stress can also improve drought stress protection. In fact this has already been demonstrated for transcription factors (in the case of AtCBF/DREB1) and for other genes such as OsCDPK7 (Saijo et al. (2000) *Plant J.* 23: 319-327), or AVP1 (a vacuolar pyrophosphatase-proton-pump; Gaxiola et al. (2001) *Proc. Natl. Acad. Sci. USA* 98: 11444-11449).

The present invention relates to methods and compositions for producing transgenic plants with modified traits, particularly traits that address agricultural and food needs. These traits, including altered sugar sensing and tolerance to abiotic stress (e.g., germination in heat or in cold conditions), and osmotic stress (e.g., tolerance to high salt concentrations or drought), may provide significant value in that the plant can then thrive in hostile environments, where, for example, high or low temperature, low water availability or high salinity may limit or prevent growth of non-transgenic plants.

We have identified polynucleotides encoding transcription factors, including G1073 (atHRC1), G1067 (AtHRC2), G2153 (AtHRC3), G2156 (AtHRC4) and their equivalents listed in the Sequence Listing, and structurally and functionally similar sequences, developed numerous transgenic plants using these polynucleotides, and have analyzed the plants for their tolerance to abiotic stresses, including those associated with heat, cold, or osmotic stresses such as drought and excessive salt. In so doing, we have identified important polynucleotide and polypeptide sequences for producing commercially valuable plants and crops as well as the methods for making them and using them. Other aspects and embodiments of the invention are described below and can be derived from the teachings of this disclosure as a whole.

SUMMARY OF THE INVENTION

The present invention pertains to recombinant polynucleotides that comprise sequences able to hybridizing under stringent conditions to the nucleotide sequences of G1073 (AtHRC1; SEQ ID NO: 1), G1067 (AtHRC2; SEQ ID NO: 3), and G2153 (AtHRC3; SEQ ID NO: 5), and their complements. These stringent conditions include 6x SSC and 65° C. These polynucleotides encode polypeptides that have the ability to regulate transcription and increase the biomass or abiotic stress tolerance of a plant.

The invention also pertains to expression vectors comprising these recombinant polynucleotides, and to cultured host plant cells that comprise these recombinant polynucleotides.

The invention is also directed to transgenic plants that comprise a recombinant polynucleotide encoding a polypeptide with an AT-hook domain. This AT-hook domain is sufficiently homologous to the AT-hook domain of G1073 (SEQ ID NO: 2) that the polypeptide is able to bind to the narrow minor groove of AT-rich regions of DNA and regulate transcription. The polypeptide also has the property of SEQ ID NO:2 in that it alters a plant's traits by regulating abiotic stress tolerance or increasing biomass in the plant. The binding of the polypeptide to the DNA being regulated ultimately confers the altered trait; plants altered in this manner may be identified by comparing a transformed plant to a non-transformed plant that does not overexpress the polypeptide. The recombinant polynucleotide sequences of the invention comprise nucleotide sequences that are capable of hybridizing over their full length to the complement of SEQ ID NO:1, SEQ ID NO: 3, SEQ ID NO: 5, SEQ ID NO: 7, SEQ ID NO: 9, SEQ ID NO: 11, SEQ ID NO: 13, SEQ ID NO: 15 or SEQ ID NO: 17 under stringent conditions comprising 6x SSC and 65° C. The polypeptides of the invention, which are encoded by these polynucleotides, include SEQ ID NO: 2, SEQ ID NO: 4, SEQ ID NO: 6, SEQ ID NO: 8, SEQ ID NO: 10, SEQ ID NO: 12, SEQ ID NO: 14, SEQ ID NO: 16 and SEQ ID NO: 18, and structurally and functionally related polypeptides.

The invention is also directed to methods for producing transgenic plants having either increased tolerance to abiotic stress or increased biomass. These method steps include first providing an expression vector that comprises: (i) a polynucleotide sequence comprising a nucleotide sequences that hybridizes its over their full length to the complement of SEQ ID NO:1, SEQ ID NO: 3, SEQ ID NO: 5, SEQ ID NO: 7, SEQ ID NO: 9, SEQ ID NO: 11, SEQ ID NO: 13, SEQ ID NO: 15 or SEQ ID NO: 17 under stringent conditions comprising 6x SSC and 65° C; and (ii) regulatory elements flanking the polynucleotide sequence, the regulatory elements being effective to control expression of the polynucleotide sequence in a target plant. The expression vector is then introduced into plant cells and the plant cells are regenerated into plants, after which the plant overexpress a polypeptide encoded by the recombinant polynucleotide. Plants with the desired altered traits (i.e., abiotic stress tolerance or increased biomass) may be identified by comparison to one or more non-transformed plants that do not overexpress the polypeptide. Plants with

desired levels of abiotic stress tolerance or increased biomass may then be selected. These method steps may further comprise crossing one of the transgenic plants with either itself or another plant, then selecting seed that develops as a result of this crossing. Progeny plants may be grown from the seed, thus producing a transgenic progeny plant having the desired altered trait of increased tolerance to abiotic stress or increased biomass.

Brief Description of the Sequence Listing and Drawings

The file of this patent contains at least one drawing executed in color. Copies of this patent with color drawing(s) will be provided by the Patent and Trademark Office upon request and payment of the necessary fee.

The Sequence Listing provides exemplary polynucleotide and polypeptide sequences of the invention. The traits associated with the use of the sequences are included in the Examples.

CD-ROM1 is a read-only memory computer-readable compact disc and contains a copy of the Sequence Listing in ASCII text format. The Sequence Listing is named "MBI0034CIP.ST25.txt" and is 153 kilobytes in size. The copies of the Sequence Listing on the CD-ROM disc are hereby incorporated by reference in their entirety.

Figure 1 shows a conservative estimate of phylogenetic relationships among the orders of flowering plants (modified from Angiosperm Phylogeny Group (1998) *Ann. Missouri Bot. Gard.* 84: 1-49). Those plants with a single cotyledon (monocots) are a monophyletic clade nested within at least two major lineages of dicots; the eudicots are further divided into rosids and asterids. *Arabidopsis* is a rosid eudicot classified within the order Brassicales; rice is a member of the monocot order Poales. Figure 1 was adapted from Daly et al. (2001) *Plant Physiol.* 127: 1328-1333.

Figure 2 shows a phylogenic dendrogram depicting phylogenetic relationships of higher plant taxa, including clades containing tomato and *Arabidopsis*; adapted from Ku et al. (2000) *Proc. Natl. Acad. Sci.* 97: 9121-9126; and Chase et al. (1993) *Ann. Missouri Bot. Gard.* 80: 528-580.

Figure 3 shows crop orthologs that were identified through BLAST analysis of proprietary and public data sources. A phylogeny tree was then generated using ClustalX based on whole protein sequences. Sequences that begin with the capital letter "G" refer to *Arabidopsis* sequences (with regard to the sequence "GID" number); "GM" refers to soy sequences, "OS" to rice sequences, and "ZM" to corn sequences. Sequences that are underlined have been shown to confer increased biomass when overexpressed. The designations G3401, OS AP004587 and OS C2099_1 all refer to the same sequence.

Figure 4 depicts the domain structure of AT-hook proteins, represented by a schematic representation of the G1073 (AtHRC1) protein. Arrows indicate potential CK2 and PKC phosphorylation

sites. A conservative DNA binding domain is located at positions 34 through 42.

In Figures 5A-5J, the alignments of the AT-hook proteins identified in Figure 3, are shown, and include *Arabidopsis* (G1073, G1067, G2153, G2156), soy G3456, G3459, G3460), and rice (G3399, G3407) sequences that have been shown to confer similar traits in plants when overexpressed. Residues that appear in boxes are conserved between these sequences, being identical or similar. Also shown are sequence alignments with other *Arabidopsis* aligned with soybean, rice and corn sequences, showing the AT-hook conserved domains (Figure 5D) and the second conserved domains spanning Figures 5E through 5G).

Figures 6A and 6B show wild-type (left) and G1073-overexpressing (right) *Arabidopsis* stem cross-sections. In the stem from the G1073-overexpressing plant, the vascular bundles are larger (containing more cells in the phloem and xylem areas) and the cells of the cortex are enlarged.

Many *Arabidopsis* plants that overexpress G1073 (Figure 7A, example on right) are larger than wild-type control plants (Figure 7A, left). This distinction also holds true for the floral organs, which, as seen in Figure 7B, are significantly larger in the G1073-overexpressing plant on the right than in that from the wild-type plant on the left.

Comparing Figures 8A and 8B, 35S::G1073 lines are seen to have increased resistance to drought related stresses. Ten of ten 35S::G1073 seedlings tested showed enhanced growth, as indicated by greater cotyledon expansion and root development, in germination assays on 150 mM NaCl. Similar results were obtained with five of ten lines on 9.4% sucrose plates (not shown).

Paralogs of G1073, including G1067, G2153 and G2156, also confer an increase in biomass when these genes are overexpressed and the plants compared with wild-type plants. G2156, for example, produces increased floral organ size (Figure 9A, overexpressors left and center) and larger plants (Figure 9B, overexpressor on left).

Figure 10 is a graph comparing silique number in control (wild type) and 35S::G1073 plants indicating how seed number is associated with the increased number of siliques per plant seen in the overexpressing lines.

As seen in Figures 11A and 11B, G1073 functions in both soybean and tomato to increase biomass. In Figure 11A, the larger plant on the right is overexpressing G1073. Tomato leaves of a number of G1073 overexpressor lines were much larger than those of wild-type tomato plants, as seen in Figure 11B by comparing the leaves of the overexpressor plant on the left and that from a wild-type plant on the right.

Figure 12A is a photograph of an *Arabidopsis* plant overexpressing the monocot gene G3399, a rice ortholog of G1073. The phenotype of increased size and mass is the same as the phenotype conferred

by *Arabidopsis* G1073 and its paralog sequences G1067, G2153 and G2157. Figure 12B similarly shows the effects of another rice ortholog, G3407, at seven days. The overexpressor on the left is approximately 50% larger than the control plant on the right.

Figure 13 shows the effects of overexpression of G3460, a soy ortholog of G1073, on plant morphology. Thirty-eight days after planting, the overexpressor on the left has significantly broader and more massive leaves than the control plant on the right. The overexpressor also demonstrates late development, a characteristic also seen when G1073 or its paralogs are overexpressed.

DETAILED DESCRIPTION OF EXEMPLARY EMBODIMENTS

The present invention relates to polynucleotides and polypeptides for modifying phenotypes of plants, particularly those associated with increased biomass and/or abiotic stress tolerance. Throughout this disclosure, various information sources are referred to and/or are specifically incorporated. The information sources include scientific journal articles, patent documents, textbooks, and World Wide Web browser-inactive page addresses, for example. While the reference to these information sources clearly indicates that they can be used by one of skill in the art, each and every one of the information sources cited herein are specifically incorporated in their entirety, whether or not a specific mention of "incorporation by reference" is noted. The contents and teachings of each and every one of the information sources can be relied on and used to make and use embodiments of the invention.

As used herein and in the appended claims, the singular forms "a", "an", and "the" include the plural reference unless the context clearly dictates otherwise. Thus, for example, a reference to "a host cell" includes a plurality of such host cells, and a reference to "a stress" is a reference to one or more stresses and equivalents thereof known to those skilled in the art, and so forth

DEFINITIONS

"Nucleic acid molecule" refers to a oligonucleotide, polynucleotide or any fragment thereof. It may be DNA or RNA of genomic or synthetic origin, double-stranded or single-stranded, and combined with carbohydrate, lipids, protein, or other materials to perform a particular activity such as transformation or form a useful composition such as a peptide nucleic acid (PNA).

"Polynucleotide" is a nucleic acid molecule comprising a plurality of polymerized nucleotides, e.g., at least about 15 consecutive polymerized nucleotides, optionally at least about 30 consecutive nucleotides, at least about 50 consecutive nucleotides. A polynucleotide may be a nucleic acid,

oligonucleotide, nucleotide, or any fragment thereof. In many instances, a polynucleotide comprises a nucleotide sequence encoding a polypeptide (or protein) or a domain or fragment thereof. Additionally, the polynucleotide may comprise a promoter, an intron, an enhancer region, a polyadenylation site, a translation initiation site, 5' or 3' untranslated regions, a reporter gene, a selectable marker, or the like. The polynucleotide can be single stranded or double stranded DNA or RNA. The polynucleotide optionally comprises modified bases or a modified backbone. The polynucleotide can be, e.g., genomic DNA or RNA, a transcript (such as an mRNA), a cDNA, a PCR product, a cloned DNA, a synthetic DNA or RNA, or the like. The polynucleotide can be combined with carbohydrate, lipids, protein, or other materials to perform a particular activity such as transformation or form a useful composition such as a peptide nucleic acid (PNA). The polynucleotide can comprise a sequence in either sense or antisense orientations. "Oligonucleotide" is substantially equivalent to the terms amplimer, primer, oligomer, element, target, and probe and is preferably single stranded.

"Gene" or "gene sequence" refers to the partial or complete coding sequence of a gene, its complement, and its 5' or 3' untranslated regions. A gene is also a functional unit of inheritance, and in physical terms is a particular segment or sequence of nucleotides along a molecule of DNA (or RNA, in the case of RNA viruses) involved in producing a polypeptide chain. The latter may be subjected to subsequent processing such as splicing and folding to obtain a functional protein or polypeptide. A gene may be isolated, partially isolated, or be found with an organism's genome. By way of example, a transcription factor gene encodes a transcription factor polypeptide, which may be functional or require processing to function as an initiator of transcription.

Operationally, genes may be defined by the cis-trans test, a genetic test that determines whether two mutations occur in the same gene and which may be used to determine the limits of the genetically active unit (Rieger et al. (1976) Glossary of Genetics and Cytogenetics: Classical and Molecular, 4th ed., Springer Verlag, Berlin). A gene generally includes regions preceding ("leaders"; upstream) and following ("trailers"; downstream) of the coding region. A gene may also include intervening, non-coding sequences, referred to as "introns", located between individual coding segments, referred to as "exons". Most genes have an associated promoter region, a regulatory sequence 5' of the transcription initiation codon (there are some genes that do not have an identifiable promoter). The function of a gene may also be regulated by enhancers, operators, and other regulatory elements.

A "recombinant polynucleotide" is a polynucleotide that is not in its native state, e.g., the polynucleotide comprises a nucleotide sequence not found in nature, or the polynucleotide is in a context other than that in which it is naturally found, e.g., separated from nucleotide sequences with which it typically is in proximity in nature, or adjacent (or contiguous with) nucleotide sequences with which it

typically is not in proximity. For example, the sequence at issue can be cloned into a vector, or otherwise recombined with one or more additional nucleic acid.

An "isolated polynucleotide" is a polynucleotide whether naturally occurring or recombinant, that is present outside the cell in which it is typically found in nature, whether purified or not. Optionally, an isolated polynucleotide is subject to one or more enrichment or purification procedures, e.g., cell lysis, extraction, centrifugation, precipitation, or the like.

A "polypeptide" is an amino acid sequence comprising a plurality of consecutive polymerized amino acid residues e.g., at least about 15 consecutive polymerized amino acid residues, optionally at least about 30 consecutive polymerized amino acid residues, at least about 50 consecutive polymerized amino acid residues. In many instances, a polypeptide comprises a polymerized amino acid residue sequence that is a transcription factor or a domain or portion or fragment thereof. Additionally, the polypeptide may comprise 1) a localization domain, 2) an activation domain, 3) a repression domain, 4) an oligomerization domain, or 5) a DNA-binding domain, or the like. The polypeptide optionally comprises modified amino acid residues, naturally occurring amino acid residues not encoded by a codon, non-naturally occurring amino acid residues.

"Protein" refers to an amino acid sequence, oligopeptide, peptide, polypeptide or portions thereof whether naturally occurring or synthetic.

"Portion", as used herein, refers to any part of a protein used for any purpose, but especially for the screening of a library of molecules which specifically bind to that portion or for the production of antibodies.

A "recombinant polypeptide" is a polypeptide produced by translation of a recombinant polynucleotide. A "synthetic polypeptide" is a polypeptide created by consecutive polymerization of isolated amino acid residues using methods well known in the art. An "isolated polypeptide," whether a naturally occurring or a recombinant polypeptide, is more enriched in (or out of) a cell than the polypeptide in its natural state in a wild-type cell, e.g., more than about 5% enriched, more than about 10% enriched, or more than about 20%, or more than about 50%, or more, enriched, i.e., alternatively denoted: 105%, 110%, 120%, 150% or more, enriched relative to wild type standardized at 100%. Such an enrichment is not the result of a natural response of a wild-type plant. Alternatively, or additionally, the isolated polypeptide is separated from other cellular components with which it is typically associated, e.g., by any of the various protein purification methods herein.

"Homology" refers to sequence similarity between a reference sequence and at least a fragment of a newly sequenced clone insert or its encoded amino acid sequence.

"Hybridization complex" refers to a complex between two nucleic acid molecules by virtue of the

formation of hydrogen bonds between purines and pyrimidines.

"Identity" or "similarity" refers to sequence similarity between two polynucleotide sequences or between two polypeptide sequences, with identity being a more strict comparison. The phrases "percent identity" and "% identity" refer to the percentage of sequence similarity found in a comparison of two or more polynucleotide sequences or two or more polypeptide sequences. "Sequence similarity" refers to the percent similarity in base pair sequence (as determined by any suitable method) between two or more polynucleotide sequences. Two or more sequences can be anywhere from 0-100% similar, or any integer value therebetween. Identity or similarity can be determined by comparing a position in each sequence that may be aligned for purposes of comparison. When a position in the compared sequence is occupied by the same nucleotide base or amino acid, then the molecules are identical at that position. A degree of similarity or identity between polynucleotide sequences is a function of the number of identical or matching nucleotides at positions shared by the polynucleotide sequences. A degree of identity of polypeptide sequences is a function of the number of identical amino acids at positions shared by the polypeptide sequences. A degree of homology or similarity of polypeptide sequences is a function of the number of amino acids at positions shared by the polypeptide sequences.

The term "amino acid consensus motif" refers to the portion or subsequence of a polypeptide sequence that is substantially conserved among the polypeptide transcription factors listed in the Sequence Listing.

"Alignment" refers to a number of nucleotide bases or amino acid residue sequences aligned by lengthwise comparison so that components in common (i.e., nucleotide bases or amino acid residues) may be visually and readily identified. The fraction or percentage of components in common is related to the homology or identity between the sequences. Alignments such as those of Figures 3, 4, or 5 may be used to identify conserved domains and relatedness within these domains. An alignment may suitably be determined by means of computer programs known in the art, such as MACVECTOR software (1999) (Accelrys, Inc., San Diego, CA).

A "conserved domain" or "conserved region" as used herein refers to a region in heterologous polynucleotide or polypeptide sequences where there is a relatively high degree of sequence identity between the distinct sequences. An "AT-hook" domain, such as is found in a member of AT-hook transcription factor family, is an example of a conserved domain. With respect to polynucleotides encoding presently disclosed transcription factors, a conserved domain is preferably at least 10 base pairs (bp) in length. A "conserved domain", with respect to presently disclosed AT-hook polypeptides refers to a domain within a transcription factor family that exhibits a higher degree of sequence homology, such as at least 62% sequence identity including conservative substitutions, and more preferably at least 65%

sequence identity, and even more preferably at least 69%, or at least about 71%, or at least about 78%, or at least about 81%, or at least about 90%, or at least about 95%, or at least about 98% amino acid residue sequence identity to the conserved domain. A fragment or domain can be referred to as outside a conserved domain, outside a consensus sequence, or outside a consensus DNA-binding site that is known to exist or that exists for a particular transcription factor class, family, or sub-family. In this case, the fragment or domain will not include the exact amino acids of a consensus sequence or consensus DNA-binding site of a transcription factor class, family or sub-family, or the exact amino acids of a particular transcription factor consensus sequence or consensus DNA-binding site. Furthermore, a particular fragment, region, or domain of a polypeptide, or a polynucleotide encoding a polypeptide, can be "outside a conserved domain" if all the amino acids of the fragment, region, or domain fall outside of a defined conserved domain(s) for a polypeptide or protein. Sequences having lesser degrees of identity but comparable biological activity are considered to be equivalents.

As one of ordinary skill in the art recognizes, conserved domains may be identified as regions or domains of identity to a specific consensus sequence (see, for example, Riechmann et al. (2000) *supra*). Thus, by using alignment methods well known in the art, the conserved domains of the plant transcription factors for the AT-hook proteins (Reeves and Beckerbauer (2001) *Biochim. Biophys. Acta* 1519: 13-29; and Reeves (2001) *Gene* 277: 63-81) may be determined.

The conserved domains for SEQ ID NO: 2, 4, 6, 8, 10, 12, 14, 16 and 18 are listed in Table 1. Also, the polypeptides of Table 1 have AT-hook and second conserved domains specifically indicated by start and stop sites. A comparison of the regions of the polypeptides in SEQ ID NO: 2, 4, 6, 8, 10, 12, 14, 16 and 18 allows one of skill in the art (see, for example, Reeves and Nisson (1995) *Biol. Chem.* 265: 8573-8582) to identify AT-hook domains or conserved domains for any of the polypeptides listed or referred to in this disclosure.

"Complementary" refers to the natural hydrogen bonding by base pairing between purines and pyrimidines. For example, the sequence A-C-G-T (5' -> 3') forms hydrogen bonds with its complements A-C-G-T (5' -> 3') or A-C-G-U (5' -> 3'). Two single-stranded molecules may be considered partially complementary, if only some of the nucleotides bond, or "completely complementary" if all of the nucleotides bond. The degree of complementarity between nucleic acid strands affects the efficiency and strength of the hybridization and amplification reactions. "Fully complementary" refers to the case where bonding occurs between every base pair and its complement in a pair of sequences, and the two sequences have the same number of nucleotides.

The terms "highly stringent" or "highly stringent condition" refer to conditions that permit hybridization of DNA strands whose sequences are highly complementary, wherein these same conditions

exclude hybridization of significantly mismatched DNAs. Polynucleotide sequences capable of hybridizing under stringent conditions with the polynucleotides of the present invention may be, for example, variants of the disclosed polynucleotide sequences, including allelic or splice variants, or sequences that encode orthologs or paralogs of presently disclosed polypeptides. Nucleic acid hybridization methods are disclosed in detail by Kashima et al. (1985) *Nature* 313:402-404, and Sambrook et al. (1989) Molecular Cloning: A Laboratory Manual, 2nd Ed., Cold Spring Harbor Laboratory, Cold Spring Harbor, N.Y. ("Sambrook"); and by Haymes et al. "Nucleic Acid Hybridization: A Practical Approach", IRL Press, Washington, D.C. (1985), which references are incorporated herein by reference.

In general, stringency is determined by the temperature, ionic strength, and concentration of denaturing agents (e.g., formamide) used in a hybridization and washing procedure (for a more detailed description of establishing and determining stringency, see below). The degree to which two nucleic acids hybridize under various conditions of stringency is correlated with the extent of their similarity. Thus, similar nucleic acid sequences from a variety of sources, such as within a plant's genome (as in the case of paralogs) or from another plant (as in the case of orthologs) that may perform similar functions can be isolated on the basis of their ability to hybridize with known transcription factor sequences. Numerous variations are possible in the conditions and means by which nucleic acid hybridization can be performed to isolate transcription factor sequences having similarity to transcription factor sequences known in the art and are not limited to those explicitly disclosed herein. Such an approach may be used to isolate polynucleotide sequences having various degrees of similarity with disclosed transcription factor sequences, such as, for example, encoded transcription factors having 62% or greater identity with the AT-hook domain of disclosed transcription factors.

Regarding the terms "paralog" and "ortholog", homologous polynucleotide sequences and homologous polypeptide sequences may be paralogs or orthologs of the claimed polynucleotide or polypeptide sequence. Orthologs and paralogs are evolutionarily related genes that have similar sequence and similar functions. Orthologs are structurally related genes in different species that are derived by a speciation event. Paralogs are structurally related genes within a single species that are derived by a duplication event. Sequences that are sufficiently similar to one another will be appreciated by those of skill in the art and may be based upon percentage identity of the complete sequences, percentage identity of a conserved domain or sequence within the complete sequence, percentage similarity to the complete sequence, percentage similarity to a conserved domain or sequence within the complete sequence, and/or an arrangement of contiguous nucleotides or peptides particular to a conserved domain or complete sequence. Sequences that are sufficiently similar to one another will also bind in a similar manner to the

same DNA binding sites of transcriptional regulatory elements using methods well known to those of skill in the art.

The term "equivalog" describes members of a set of homologous proteins that are conserved with respect to function since their last common ancestor. Related proteins are grouped into equivalog families, and otherwise into protein families with other hierarchically defined homology types. This definition is provided at the Institute for Genomic Research (TIGR) world wide web (www) website, " tigr.org " under the heading "Terms associated with TIGRFAMs".

The term "variant", as used herein, may refer to polynucleotides or polypeptides, that differ from the presently disclosed polynucleotides or polypeptides, respectively, in sequence from each other, and as set forth below.

With regard to polynucleotide variants, differences between presently disclosed polynucleotides and polynucleotide variants are limited so that the nucleotide sequences of the former and the latter are closely similar overall and, in many regions, identical. Due to the degeneracy of the genetic code, differences between the former and latter nucleotide sequences may be silent (i.e., the amino acids encoded by the polynucleotide are the same, and the variant polynucleotide sequence encodes the same amino acid sequence as the presently disclosed polynucleotide. Variant nucleotide sequences may encode different amino acid sequences, in which case such nucleotide differences will result in amino acid substitutions, additions, deletions, insertions, truncations or fusions with respect to the similar disclosed polynucleotide sequences. These variations result in polynucleotide variants encoding polypeptides that share at least one functional characteristic. The degeneracy of the genetic code also dictates that many different variant polynucleotides can encode identical and/or substantially similar polypeptides in addition to those sequences illustrated in the Sequence Listing.

Also within the scope of the invention is a variant of a transcription factor nucleic acid listed in the Sequence Listing, that is, one having a sequence that differs from the one of the polynucleotide sequences in the Sequence Listing, or a complementary sequence, that encodes a functionally equivalent polypeptide (i.e., a polypeptide having some degree of equivalent or similar biological activity) but differs in sequence from the sequence in the Sequence Listing, due to degeneracy in the genetic code. Included within this definition are polymorphisms that may or may not be readily detectable using a particular oligonucleotide probe of the polynucleotide encoding polypeptide, and improper or unexpected hybridization to allelic variants, with a locus other than the normal chromosomal locus for the polynucleotide sequence encoding polypeptide.

"Allelic variant" or "polynucleotide allelic variant" refers to any of two or more alternative forms of a gene occupying the same chromosomal locus. Allelic variation arises naturally through mutation, and

may result in phenotypic polymorphism within populations. Gene mutations may be "silent" or may encode polypeptides having altered amino acid sequence. "Allelic variant" and "polypeptide allelic variant" may also be used with respect to polypeptides, and in this case the term refer to a polypeptide encoded by an allelic variant of a gene.

5 "Splice variant" or "polynucleotide splice variant" as used herein refers to alternative forms of RNA transcribed from a gene. Splice variation naturally occurs as a result of alternative sites being spliced within a single transcribed RNA molecule or between separately transcribed RNA molecules, and may result in several different forms of mRNA transcribed from the same gene. This, splice variants may encode polypeptides having different amino acid sequences, which may or may not have similar functions
10 in the organism. "Splice variant" or "polypeptide splice variant" may also refer to a polypeptide encoded by a splice variant of a transcribed mRNA.

As used herein, "polynucleotide variants" may also refer to polynucleotide sequences that encode paralogs and orthologs of the presently disclosed polypeptide sequences. "Polypeptide variants" may refer to polypeptide sequences that are paralogs and orthologs of the presently disclosed polypeptide sequences.

15 Differences between presently disclosed polypeptides and polypeptide variants are limited so that the sequences of the former and the latter are closely similar overall and, in many regions, identical. Presently disclosed polypeptide sequences and similar polypeptide variants may differ in amino acid sequence by one or more substitutions, additions, deletions, fusions and truncations, which may be present in any combination. These differences may produce silent changes and result in a functionally equivalent
20 transcription factor. Thus, it will be readily appreciated by those of skill in the art, that any of a variety of polynucleotide sequences is capable of encoding the transcription factors and transcription factor homolog polypeptides of the invention. A polypeptide sequence variant may have "conservative" changes, wherein a substituted amino acid has similar structural or chemical properties. Deliberate amino acid substitutions may thus be made on the basis of similarity in polarity, charge, solubility, hydrophobicity, hydrophilicity,
25 and/or the amphipathic nature of the residues, as long as the functional or biological activity of the transcription factor is retained. For example, negatively charged amino acids may include aspartic acid and glutamic acid, positively charged amino acids may include lysine and arginine, and amino acids with uncharged polar head groups having similar hydrophilicity values may include leucine, isoleucine, and valine; glycine and alanine; asparagine and glutamine; serine and threonine; and phenylalanine and
30 tyrosine (for more detail on conservative substitutions, see Table 2). More rarely, a variant may have "non-conservative" changes, e.g., replacement of a glycine with a tryptophan. Similar minor variations may also include amino acid deletions or insertions, or both. Related polypeptides may comprise, for example, additions and/or deletions of one or more N-linked or O-linked glycosylation sites, or an addition and/or a

deletion of one or more cysteine residues. Guidance in determining which and how many amino acid residues may be substituted, inserted or deleted without abolishing functional or biological activity may be found using computer programs well known in the art, for example, DNASTAR software (see USPN 5,840,544).

5 "Ligand" refers to any molecule, agent, or compound that will bind specifically to a complementary site on a nucleic acid molecule or protein. Such ligands stabilize or modulate the activity of nucleic acid molecules or proteins of the invention and may be composed of at least one of the following: inorganic and organic substances including nucleic acids, proteins, carbohydrates, fats, and lipids.

10 "Modulates" refers to a change in activity (biological, chemical, or immunological) or lifespan resulting from specific binding between a molecule and either a nucleic acid molecule or a protein.

The term "plant" includes whole plants, shoot vegetative organs/structures (for example, leaves, stems and tubers), roots, flowers and floral organs/structures (for example, bracts, sepals, petals, stamens, carpels, anthers and ovules), seed (including embryo, endosperm, and seed coat) and fruit (the mature
15 ovary), plant tissue (for example, vascular tissue, ground tissue, and the like) and cells for example, guard cells, egg cells, and the like), and progeny of same. The class of plants that can be used in the method of the invention is generally as broad as the class of higher and lower plants amenable to transformation techniques, including angiosperms (monocotyledonous and dicotyledonous plants), gymnosperms, ferns, horsetails, psilophytes, lycophytes, bryophytes, and multicellular algae. (See for example, Figure 1,
20 adapted from Daly et al. (2001) *Plant Physiol.* 127: 1328-1333; Figure 2, adapted from Ku et al. (2000) *Proc. Natl. Acad. Sci.* 97: 9121-9126; and see also Tudge in The Variety of Life, Oxford University Press, New York, NY (2000) pp. 547-606).

A "transgenic plant" refers to a plant that contains genetic material not found in a wild-type plant of the same species, variety or cultivar. The genetic material may include a transgene, an insertional
25 mutagenesis event (such as by transposon or T-DNA insertional mutagenesis), an activation tagging sequence, a mutated sequence, a homologous recombination event or a sequence modified by chimera-plasty. Typically, the foreign genetic material has been introduced into the plant by human manipulation, but any method can be used as one of skill in the art recognizes.

A transgenic plant may contain an expression vector or cassette. The expression cassette typically
30 comprises a polypeptide-encoding sequence operably linked (i.e., under regulatory control of) to appropriate inducible or constitutive regulatory sequences that allow for the expression of polypeptide. The expression cassette can be introduced into a plant by transformation or by breeding after transformation of a parent plant. A plant refers to a whole plant as well as to a plant part, such as seed,

fruit, leaf, or root, plant tissue, plant cells or any other plant material, e.g., a plant explant, as well as to progeny thereof, and to *in vitro* systems that mimic biochemical or cellular components or processes in a cell.

"Control plant" refers to a plant that serves as a standard of comparison for testing the results of a treatment or genetic alteration, or the degree of altered expression of a gene or gene product. Examples of control plants include plants that are untreated, or genetically unaltered (i.e., wild-type).

"Wild type", as used herein, refers to a cell, tissue or plant that has not been genetically modified to knock out or overexpress one or more of the presently disclosed transcription factors. Wild-type cells, tissue or plants may be used as controls to compare levels of expression and the extent and nature of trait modification with cells, tissue or plants in which transcription factor expression is altered or ectopically expressed, e.g., in that it has been knocked out or overexpressed.

"Fragment", with respect to a polynucleotide, refers to a clone or any part of a polynucleotide molecule that retains a usable, functional characteristic. Useful fragments include oligonucleotides and polynucleotides that may be used in hybridization or amplification technologies or in the regulation of replication, transcription or translation. A polynucleotide fragment" refers to any subsequence of a polynucleotide, typically, of at least about 9 consecutive nucleotides, preferably at least about 30 nucleotides, more preferably at least about 50 nucleotides, of any of the sequences provided herein. Exemplary polynucleotide fragments are the first sixty consecutive nucleotides of the transcription factor polynucleotides listed in the Sequence Listing. Exemplary fragments also include fragments that comprise a region that encodes an AT-hook domain of a transcription factor. Exemplary fragments also include fragments that comprise a conserved domain of a transcription factor. Exemplary fragments include fragments that comprise an AT-hook or second conserved domain of an AT-hook transcription factor, for example, amino acid residues 34-42 and 78-175 of G1073 (AtHRC1; SEQ ID NO: 2), as noted in Table 1.

Fragments may also include subsequences of polypeptides and protein molecules, or a subsequence of the polypeptide. Fragments may have uses in that they may have antigenic potential. In some cases, the fragment or domain is a subsequence of the polypeptide which performs at least one biological function of the intact polypeptide in substantially the same manner, or to a similar extent, as does the intact polypeptide. For example, a polypeptide fragment can comprise a recognizable structural motif or functional domain such as a DNA-binding site or domain that binds to a DNA promoter region, an activation domain, or a domain for protein-protein interactions, and may initiate transcription. Fragments can vary in size from as few as 3 amino acid residues to the full length of the intact polypeptide, but are preferably at least about 30 amino acid residues in length and more preferably at least about 60 amino acid residues in length.

The invention also encompasses production of DNA sequences that encode transcription factors and transcription factor derivatives, or fragments thereof, entirely by synthetic chemistry. After production, the synthetic sequence may be inserted into any of the many available expression vectors and cell systems using reagents well known in the art. Moreover, synthetic chemistry may be used to introduce mutations into a sequence encoding transcription factors or any fragment thereof.

"Derivative" refers to the chemical modification of a nucleic acid molecule or amino acid sequence. Chemical modifications can include replacement of hydrogen by an alkyl, acyl, or amino group or glycosylation, pegylation, or any similar process that retains or enhances biological activity or lifespan of the molecule or sequence.

A "trait" refers to a physiological, morphological, biochemical, or physical characteristic of a plant or particular plant material or cell. In some instances, this characteristic is visible to the human eye, such as seed or plant size, or can be measured by biochemical techniques, such as detecting the protein, starch, or oil content of seed or leaves, or by observation of a metabolic or physiological process, e.g. by measuring tolerance to water deprivation or particular salt or sugar concentrations, or by the observation of the expression level of a gene or genes, e.g., by employing Northern analysis, RT-PCR, microarray gene expression assays, or reporter gene expression systems, or by agricultural observations such as osmotic stress tolerance or yield. Any technique can be used to measure the amount of, comparative level of, or difference in any selected chemical compound or macromolecule in the transgenic plants, however.

"Trait modification" refers to a detectable difference in a characteristic in a plant ectopically expressing a polynucleotide or polypeptide of the present invention relative to a plant not doing so, such as a wild-type plant. In some cases, the trait modification can be evaluated quantitatively. For example, the trait modification can entail at least about a 2% increase or decrease in an observed trait (difference), at least a 5% difference, at least about a 10% difference, at least about a 20% difference, at least about a 30%, at least about a 50%, at least about a 70%, or at least about a 100%, or an even greater difference compared with a wild-type plant. It is known that there can be a natural variation in the modified trait. Therefore, the trait modification observed entails a change of the normal distribution of the trait in the plants compared with the distribution observed in wild-type plants.

The term "transcript profile" refers to the expression levels of a set of genes in a cell in a particular state, particularly by comparison with the expression levels of that same set of genes in a cell of the same type in a reference state. For example, the transcript profile of a particular transcription factor in a suspension cell is the expression levels of a set of genes in a cell knocking out or overexpressing that transcription factor compared with the expression levels of that same set of genes in a suspension cell that has normal levels of that transcription factor. The transcript profile can be presented as a list of those

genes whose expression level is significantly different between the two treatments, and the difference ratios. Differences and similarities between expression levels may also be evaluated and calculated using statistical and clustering methods.

“Ectopic expression or altered expression” in reference to a polynucleotide indicates that the pattern of expression in, e.g., a transgenic plant or plant tissue, is different from the expression pattern in a wild-type plant or a reference plant of the same species. The pattern of expression may also be compared with a reference expression pattern in a wild-type plant of the same species. For example, the polynucleotide or polypeptide is expressed in a cell or tissue type other than a cell or tissue type in which the sequence is expressed in the wild-type plant, or by expression at a time other than at the time the sequence is expressed in the wild-type plant, or by a response to different inducible agents, such as hormones or environmental signals, or at different expression levels (either higher or lower) compared with those found in a wild-type plant. The term also refers to altered expression patterns that are produced by lowering the levels of expression to below the detection level or completely abolishing expression. The resulting expression pattern can be transient or stable, constitutive or inducible. In reference to a polypeptide, the term "ectopic expression or altered expression" further may relate to altered activity levels resulting from the interactions of the polypeptides with exogenous or endogenous modulators or from interactions with factors or as a result of the chemical modification of the polypeptides.

The term “overexpression” as used herein refers to a greater expression level of a gene in a plant, plant cell or plant tissue, compared to expression in a wild-type plant, cell or tissue, at any developmental or temporal stage for the gene. Overexpression can occur when, for example, the genes encoding one or more transcription factors are under the control of a strong expression signal, such as one of the promoters described herein (e.g., the cauliflower mosaic virus 35S transcription initiation region). Overexpression may occur throughout a plant or in specific tissues of the plant, depending on the promoter used, as described below.

Overexpression may take place in plant cells normally lacking expression of polypeptides functionally equivalent or identical to the present transcription factors. Overexpression may also occur in plant cells where endogenous expression of the present transcription factors or functionally equivalent molecules normally occurs, but such normal expression is at a lower level. Overexpression thus results in a greater than normal production, or “overproduction” of the transcription factor in the plant, cell or tissue.

The term “transcription regulating region” refers to a DNA regulatory sequence that regulates expression of one or more genes in a plant when a transcription factor having one or more specific binding domains binds to the DNA regulatory sequence. Transcription factors of the present invention possess an

AP2 domain, a B3 domain, or both of these binding domains. The AP2 domain of the transcription factor binds to a transcription regulating region comprising the motif CAACA, and the B3 domain of the same transcription factor binds to a transcription regulating region comprising the motif CACCTG. The transcription factors of the invention also comprise an amino acid subsequence that forms a transcription activation domain that regulates expression of one or more abiotic stress tolerance genes in a plant when the transcription factor binds to the regulating region.

The term "phase change" refers to a plant's progression from embryo to adult, and, by some definitions, the transition wherein flowering plants gain reproductive competency. It is believed that phase change occurs either after a certain number of cell divisions in the shoot apex of a developing plant, or when the shoot apex achieves a particular distance from the roots. Thus, altering the timing of phase changes may affect a plant's size, which, in turn, may affect yield and biomass.

A "sample" with respect to a material containing nucleic acid molecules may comprise a bodily fluid; an extract from a cell, chromosome, organelle, or membrane isolated from a cell; genomic DNA, RNA, or cDNA in solution or bound to a substrate; a cell; a tissue; a tissue print; a forensic sample; and the like. In this context "substrate" refers to any rigid or semi-rigid support to which nucleic acid molecules or proteins are bound and includes membranes, filters, chips, slides, wafers, fibers, magnetic or nonmagnetic beads, gels, capillaries or other tubing, plates, polymers, and microparticles with a variety of surface forms including wells, trenches, pins, channels and pores. A substrate may also refer to a reactant in a chemical or biological reaction, or a substance acted upon (e.g., by an enzyme).

"Substantially purified" refers to nucleic acid molecules or proteins that are removed from their natural environment and are isolated or separated, and are at least about 60% free, preferably about 75% free, and most preferably about 90% free, from other components with which they are naturally associated.

DETAILED DESCRIPTION

Transcription Factors Modify Expression of Endogenous Genes

A transcription factor may include, but is not limited to, any polypeptide that can activate or repress transcription of a single gene or a number of genes. As one of ordinary skill in the art recognizes, transcription factors can be identified by the presence of a region or domain of structural similarity or identity to a specific consensus sequence or the presence of a specific consensus DNA-binding site or DNA-binding site motif (see, for example, Riechmann et al. (2000) *Science* 290: 2105-2110). The plant transcription factors may belong to the AT-hook transcription factor family (Reeves and Beckerbauer (2001) *Biochim. Biophys. Acta* 1519: 13-29; and Reeves (2001) *Gene* 277: 63-81).

Generally, the transcription factors encoded by the present sequences are involved in cell

differentiation and proliferation and the regulation of growth. Accordingly, one skilled in the art would recognize that by expressing the present sequences in a plant, one may change the expression of autologous genes or induce the expression of introduced genes. By affecting the expression of similar autologous sequences in a plant that have the biological activity of the present sequences, or by

5 introducing the present sequences into a plant, one may alter a plant's phenotype to one with improved traits related to osmotic stresses. The sequences of the invention may also be used to transform a plant and introduce desirable traits not found in the wild-type cultivar or strain. Plants may then be selected for those that produce the most desirable degree of over- or under-expression of target genes of interest and coincident trait improvement.

10 The sequences of the present invention may be from any species, particularly plant species, in a naturally occurring form or from any source whether natural, synthetic, semi-synthetic or recombinant. The sequences of the invention may also include fragments of the present amino acid sequences. Where "amino acid sequence" is recited to refer to an amino acid sequence of a naturally occurring protein molecule, "amino acid sequence" and like terms are not meant to limit the amino acid sequence to the

15 complete native amino acid sequence associated with the recited protein molecule.

In addition to methods for modifying a plant phenotype by employing one or more polynucleotides and polypeptides of the invention described herein, the polynucleotides and polypeptides of the invention have a variety of additional uses. These uses include their use in the recombinant production (i.e., expression) of proteins; as regulators of plant gene expression, as diagnostic probes for

20 the presence of complementary or partially complementary nucleic acids (including for detection of natural coding nucleic acids); as substrates for further reactions, e.g., mutation reactions, PCR reactions, or the like; as substrates for cloning e.g., including digestion or ligation reactions; and for identifying exogenous or endogenous modulators of the transcription factors. In many instances, a polynucleotide comprises a nucleotide sequence encoding a polypeptide (or protein) or a domain or fragment thereof. Additionally,

25 the polynucleotide may comprise a promoter, an intron, an enhancer region, a polyadenylation site, a translation initiation site, 5' or 3' untranslated regions, a reporter gene, a selectable marker, or the like. The polynucleotide can be single stranded or double stranded DNA or RNA. The polynucleotide optionally comprises modified bases or a modified backbone. The polynucleotide can be, e.g., genomic DNA or RNA, a transcript (such as an mRNA), a cDNA, a PCR product, a cloned DNA, a synthetic DNA

30 or RNA, or the like. The polynucleotide can comprise a sequence in either sense or antisense orientations.

Expression of genes that encode transcription factors that modify expression of endogenous genes, polynucleotides, and proteins are well known in the art. In addition, transgenic plants comprising isolated polynucleotides encoding transcription factors may also modify expression of endogenous genes,

polynucleotides, and proteins. Examples include Peng et al. (1997, *Genes Development* 11: 3194-3205) and Peng et al. (1999, *Nature*, 400: 256-261). In addition, many others have demonstrated that an *Arabidopsis* transcription factor expressed in an exogenous plant species elicits the same or very similar phenotypic response. See, for example, Fu et al. (2001, *Plant Cell* 13: 1791-1802); Nandi et al. (2000, *Curr. Biol.* 10: 215-218); Coupland (1995, *Nature* 377: 482-483); and Weigel and Nilsson (1995, *Nature* 377: 482-500).

In another example, Mandel et al. (1992, *Cell* 71:133-143) and Suzuki et al.(2001, *Plant J.* 28: 409-418) teach that a transcription factor expressed in another plant species elicits the same or very similar phenotypic response of the endogenous sequence, as often predicted in earlier studies of *Arabidopsis* transcription factors in *Arabidopsis* (see Mandel et al. 1992, *supra*; Suzuki et al. 2001, *supra*).

Other examples include Müller et al. (2001, *Plant J.* 28: 169-179); Kim et al. (2001, *Plant J.* 25: 247-259); Kyoizuka and Shimamoto (2002, *Plant Cell Physiol.* 43: 130-135); Boss and Thomas (2002, *Nature*, 416: 847-850); He et al. (2000, *Transgenic Res.* 9: 223-227); and Robson et al. (2001, *Plant J.* 28: 619-631).

In yet another example, Gilmour et al. (1998, *Plant J.* 16: 433-442) teach an *Arabidopsis* AP2 transcription factor, CBF1 (SEQ ID NO: 70), which, when overexpressed in transgenic plants, increases plant freezing tolerance. Jaglo et al. (2001, *Plant Physiol.* 127: 910-917) further identified sequences in *Brassica napus* which encode CBF-like genes and that transcripts for these genes accumulated rapidly in response to low temperature. Transcripts encoding CBF-like proteins were also found to accumulate rapidly in response to low temperature in wheat, as well as in tomato. An alignment of the CBF proteins from *Arabidopsis*, *B. napus*, wheat, rye, and tomato revealed the presence of conserved consecutive amino acid residues, PKK/RPAGR_xKFxETRHP and DSAWR, that bracket the AP2/EREBP DNA binding domains of the proteins and distinguish them from other members of the AP2/EREBP protein family. (See Jaglo et al. *supra*.)

Transcription factors mediate cellular responses and control traits through altered expression of genes containing cis-acting nucleotide sequences that are targets of the introduced transcription factor. It is well appreciated in the Art that the effect of a transcription factor on cellular responses or a cellular trait is determined by the particular genes whose expression is either directly or indirectly (e.g., by a cascade of transcription factor binding events and transcriptional changes) altered by transcription factor binding. In a global analysis of transcription comparing a standard condition with one in which a transcription factor is overexpressed, the resulting transcript profile associated with transcription factor overexpression is related to the trait or cellular process controlled by that transcription factor. For example, the PAP2 gene (and other genes in the MYB family) have been shown to control anthocyanin biosynthesis through regulation

of the expression of genes known to be involved in the anthocyanin biosynthetic pathway (Bruce et al. (2000) *Plant Cell* 12: 65-79; and Borevitz et al. (2000) *Plant Cell* 12: 2383-2393). Further, global transcript profiles have been used successfully as diagnostic tools for specific cellular states (e.g., cancerous vs. non-cancerous; Bhattacharjee et al. (2001) *Proc. Natl. Acad. Sci. USA* 98: 13790-13795; and Xu et al. (2001) *Proc Natl Acad Sci, USA* 98: 15089-15094). Consequently, it is evident to one skilled in the art that similarity of transcript profile upon overexpression of different transcription factors would indicate similarity of transcription factor function.

The AT-hook transcription factor family

In higher organisms, genomic DNA is assembled into multilevel complexes with a range of DNA-binding proteins, including the well-known histones and non-histone proteins such as the high mobility group (HMG) proteins. HMG proteins are classified into different groups based on their DNA-binding motifs, and one such group is the HMG-I(Y) subgroup (recently renamed as HMGA). Proteins in this group have been shown to bind to the minor groove of DNA via a conserved nine amino acid peptide (KRPRGRPCK) called the AT-hook motif (Reeves and Nisson (1995) *Biol. Chem.* 265: 8573-8582). At the center of this AT-hook motif is a short, strongly conserved tripeptide of glycine-arginine-proline (GRP). This simple AT-hook motif can be present in a variable number of copies (1-15) in a given AT-hook protein. For example, the mammalian HMGA1 protein has three copies of this motif. The mammalian HMGA proteins participate in a wide variety of nuclear processes ranging from chromosome and chromatin remodeling, to acting as architectural transcription factors that regulate the expression of numerous genes in vivo. As a result, these proteins influence a diverse array of cellular processes including growth, proliferation, differentiation and death through the protein-DNA and protein-protein interactions (for reviews, see Reeves and Beckerbauer (2001) *Biochim. Biophys. Acta* 1519: 13-29; and Reeves (2001) *Gene* 277: 63-81). It has been shown that HMGA proteins specifically interact with a large number of other proteins, most of which are transcription factors (Reeves (2001) *supra*). They are also subject to many types of post-translational modification. One example is phosphorylation, which markedly influences their ability to interact with DNA substrates, other proteins, and chromatin (Onate et al. (1994) *Mol. Cell Biol.* 14: 3376-3391; Falvo et al. (1995) *Cell* 83: 1101-1111; Reeves and Nissen (1995) *supra*; Huth et al. (1997) *Nat. Struct. Biol.* 4, 657-665; and Girard et al. (1998) *EMBO J.* 17: 2079-2085).

In plants, a protein with AT-hook DNA-binding motifs was identified in oat (Nieto-Sotelo and Quail (1994) *Biochem. Soc. Symp.* 60, 265-275). This protein binds to the PE1 region in the oat phytochrome A3 gene promoter, and may be involved in positive regulation of *PHYA3* gene expression (Nieto-Sotelo and Quail (1994) *supra*). DNA-binding proteins containing AT-hook domains have also

been identified in a variety of plant species, including rice, pea and *Arabidopsis* (Meijer et al. (1996) *Plant Mol. Biol.* 31: 607-618; and Gupta et al (1997a) *Plant Mol. Biol.* 35, 987-992). The rice AT-hook genes are predominantly expressed in young and meristematic tissues, suggesting that AT-hook proteins may affect the expression of genes that determine the differentiation status of cells. The pea AT-hook gene is expressed in all organs including roots, stems, leaves, flowers, tendrils and developing seeds (Gupta et al. (1997a) *supra*). Northern blot analysis revealed that an *Arabidopsis* AT-hook gene was expressed in all organs with the highest expression in flowers and developing siliques (Gupta et al. (1997b) *Plant Mol. Biol.* 34: 529-536).

To date, relatively little public data is available regarding the function of AT-hook proteins. However, an activation tagged mutant for an *Arabidopsis* AT-hook gene (corresponding to G1067, SEQ ID NO: 4) has been identified by Weigel et al. ((2000) *Plant Physiol.* 122, 1003-1013). In this G1067 activation line, delayed flowering was observed, and leaves were wavy, dark green, larger, and rounder than in wild type. Moreover, both leaf petioles and stem internodes were shorter in this line than wild type. Such complex phenotypes suggest that the gene influences a wide range of developmental processes.

Recently, it has also been shown that expression of a maize AT-hook protein in yeast cells produces better growth on a medium containing high nickel concentrations. Such an effect suggests that the protein might have influence chromatin structure, and thereby restrict nickel ion accessibility to DNA (Forzani et al. (2001)). *J. Biol. Chem.* 276, 16731-16738).

Novel AT-hook transcription factor genes and binding motifs in *Arabidopsis* and other diverse species

To date, we have identified at least thirty-four *Arabidopsis* genes that code for proteins with AT-hook DNA-binding motifs. Of these, there are twenty-two genes encoding a single AT-hook DNA-binding motif; eight genes encoding two AT-hook DNA-binding motifs; three genes (G280, G1367 and G2787, SEQ ID NOs: 55, 57 and 59, respectively) encoding four AT-hook DNA-binding motifs and a single gene (G3045, SEQ ID NO: 61) encoding three AT-hook DNA-binding motifs.

G1073 (AtHRC1), for example, contains a single typical AT-hook DNA-binding motif (RRPRGRPAG) corresponding to positions 34 to 42 within the protein. A highly conserved 129 amino acid residue domain with unknown function (henceforth referred to as the "second conserved domain") can be identified in the single AT-hook domain subgroup. Following this region, a potential acidic domain spans from position 172 to 190. Additionally, analysis of the protein using PROSITE reveals three potential protein kinase C phosphorylation sites at Ser32, Thr83 and Thr102, and three potential casein kinase II phosphorylation sites at Ser6, Ser70 and Ser247 (Figure 3). Compared to many other AT-hook proteins, the G1073 protein contains a shorter N-terminus (Figures 5A-5C).

Members of the G1073 clade are structurally distinct from other AT-hook-related proteins (as may be seen in Figures 5E-5G, comparing G1068 and above sequences near the top of the alignment, and BAB64709 and G3462 near the bottom of the alignment, with this clade in the middle of the alignment.

5 Table 1 shows the polypeptides identified by: (a) polypeptide SEQ ID NO.; (b) Gene ID (GID) No.; (c) the conserved domain coordinates for the AT-hook and second conserved domain in amino acid residue coordinates and, for G1073, G1067 and G2153, polynucleotide base coordinates encoding the conserved domains; (d) AT-hook sequences of the respective polypeptides; (e) the identity in percentage terms to the AT-hook domain of G1073; (f) second conserved domain sequences of the respective
10 polypeptides; and (g) the identity in percentage terms to the second conserved domain of G1073.

Table 1. Gene families and binding domains

SEQ ID NO:	GID No.	AT-hook and Second Conserved Domains in AA Coordinates and Base Coordinates	First domain	% ID to First Domain of G1073	Second Conserved Domain	% ID to Second Conserved Domain of G1073
2	G1073 AtHRC1	Polypeptide coordinates: 34-42; 78-175 Polynucleotide coordinates: 161-187; 293-586	RRPRGRPAG	100%	VSTYATRRGCGVCIISGT GAVTNVTIRQPAAPAGG GVITLHGRFDILSLTGTA LPPPAPPAGGLTVYLA GGQGQVVGGNVAGSLI ASGPVVLMAASF	100%
4	G1067 AtHRC2	Polypeptide coordinates: 86-94, 130-235 Polynucleotide coordinates: 691-717; 823-1137	KRPRGRPPG	78%	VSTYARRRGRGVSVLG GNGTVSNVTLRQPVTGP NGGGVSGGGGVVTLHG RFEILSLTGTVLPPPAPP GAGGLSIFLAGGQGQVV GGSVVAPLIASAPVILM AASF	69%
6	G2153 AtHRC3	Polypeptide coordinates: 80-88, 124-227 Polynucleotide coordinates: 480-506; 612-923	RRPRGRPAG	89%	LATFARRRQRGICILSGN GTVANVTLRQPSTAAVA AAPGGA AVLALQGRFEI LSLTGSFLPGPAPPGSTG LTIYLAGGQGVVGGSV VGPLMAAGPVMLIAATF	62%
8	G2156 AtHRC4	Polypeptide coordinates: 72-80, 116-220	KRPRGRPPG	78%	VTTYARRRGRGVSILSG NGTVANVSLRQPATTAA HGANGGTGGVVALHGR FEILSLTGTVLPPPAPPGS GGLSIFLSGVQGVIGG NVVAPLVASGPVILMAA SF	65%
10	G3399	Polypeptide coordinates: 99-107, 143-240:	RRPRGRPPG	78%	VAEYARRRGRGVCVLS GGGAVVNVALRQPGAS PPGSMVATLRGRFEILSL TGTVLPPPAPPGASGLT	71%

					VFLSGGQGGVIGGSVVG PLVAAGPVVLMMAAS	
12	G3407	Polypeptide coordinates: 63-71, 106-208	RRPRGRPPG	78%	LTAYARRRQRGVCVLSA AGTVANVTLRQPQSAQP GPASPAVATLHGRFEILS LAGSFLPPPAPPGATSLA AFLAGGQGGVVGGSVA GALIAAGPVVVVAASF	63%
14	G3456	Polypeptide coordinates: 62-70, 106-201	RRPRGRPPG	78%	VAQFARRRQRGVSILSG SGTVVNVNLRQPTAPGA VMALHGRFDILSLTGSF LPGSPPGATGLTIYLAG GQGQIVGGEVVGPLVA AGPVLVMAATF	65%
16	G3459	Polypeptide coordinates: 76-84, 121-216	RRPRGRPPG	89%	VTAYARRRQRGICVLSG SGTVTNVSLRQPAAAGA VVTLHGRFEILSLSGSFL PPPAPPGATSLTIYLAGG QGQVVGGNVIGELTAA GPVIVIAASF	68%
18	G3460	Polypeptide coordinates: 74-82, 118-213	RRPRGRPSG	89%	VTAYARRRQRGICVLSG SGTVTNVSLRQPAAAGA VRLHGRFEILSLSGSFL PPPAPPGATSLTIYLAGG QGQVVGGNVVIGELTAA GPVIVIAASF	67%

The transcription factors of the invention each possess an AT-hook domain comprising two conserved domains, and include paralogs and orthologs of G1073 found by BLAST analysis, as described below. As shown in Table 1, the AT-hook domains of G1073 and related sequences are at least 78% identical to the At-Hook domains of G1073 and at least 62% identical to the second conserved domain found in G1073. These transcription factors rely on the binding specificity of their AT-hook domains, all have been shown to similar or identical functions in plants by increasing the size and biomass of a plant.

Polypeptides and Polynucleotides of the Invention

The present invention provides, among other things, transcription factors (TFs), and transcription factor homolog polypeptides, and isolated or recombinant polynucleotides encoding the polypeptides, or novel sequence variant polypeptides or polynucleotides encoding novel variants of transcription factors derived from the specific sequences provided in the Sequence Listing. Also provided are methods for modifying a plant's biomass by modifying the size or number of leaves or seed of a plant by controlling a number of cellular processes, and for increasing a plant's tolerance to abiotic stresses. This is achieved by altering the expression of critical regulatory molecules that may be conserved between diverse plant species; related conserved regulatory molecules may be originally discovered in a model system such as

Arabidopsis and homologous, functional molecules then discovered in other plant species.

The polypeptide and polynucleotide sequences of G1067 were previously identified in U.S. Provisional Patent Application 60/135,134, filed May 20, 1999. The polypeptide and polynucleotide sequences of G1073 were previously identified in U.S. Provisional Patent Application 60/125,814, filed March 23, 1999. The function of G1073 in increasing biomass was disclosed in US Provisional Application No. 60/227,439, filed August 22, 2000, and the utility for increased drought tolerance observed in 35S::G1073 transgenic lines was disclosed in US Non-Provisional Application No. 10/374,780, filed February 25, 2003. The polypeptide and polynucleotide sequences of G2153 and G2156 were previously identified in U.S. Provisional Patent Application No. 60/338,692, filed December 11, 2001, and in U.S. Non-provisional Patent Applications 10/225,066 and 10/225,068, both of which were filed August 9, 2002. The altered sugar sensing and osmotic stress tolerance phenotype conferred by G2153 overexpression was disclosed in these filings. At the time each of the above applications were filed, these sequences were identified as encoding or being transcription factors, which were defined as polypeptides having the ability to effect transcription of a target gene. It is noted that sequences that have gene-regulating activity have been determined to have specific and substantial utility by the U.S. Patent and Trademark Office (*Federal Register* (2001) 66(4): 1095).

Exemplary polynucleotides encoding the polypeptides of the invention were identified in the *Arabidopsis thaliana* GenBank database using publicly available sequence analysis programs and parameters. Sequences initially identified were then further characterized to identify sequences comprising specified sequence strings corresponding to sequence motifs present in families of known transcription factors. In addition, further exemplary polynucleotides encoding the polypeptides of the invention were identified in the plant GenBank database using publicly available sequence analysis programs and parameters. Sequences initially identified were then further characterized to identify sequences comprising specified sequence strings corresponding to sequence motifs present in families of known transcription factors. Polynucleotide sequences meeting such criteria were confirmed as transcription factors.

Additional polynucleotides of the invention were identified by screening *Arabidopsis thaliana* and/or other plant cDNA libraries with probes corresponding to known transcription factors under low stringency hybridization conditions. Additional sequences, including full length coding sequences were subsequently recovered by the rapid amplification of cDNA ends (RACE) procedure, using a commercially available kit according to the manufacturer's instructions. Where necessary, multiple rounds of RACE are performed to isolate 5' and 3' ends. The full-length cDNA was then recovered by a routine end-to-end polymerase chain reaction (PCR) using primers specific to the isolated 5' and 3' ends. Exemplary sequences are provided in the Sequence Listing.

The polynucleotides of the invention can be or have been ectopically expressed in overexpressor or knockout plants and the changes in the characteristic(s) or trait(s) of the plants observed. Therefore, the polynucleotides and polypeptides can be employed to improve the characteristics of plants.

The polynucleotides of the invention can be or were ectopically expressed in overexpressor plant cells and the changes in the expression levels of a number of genes, polynucleotides, and/or proteins of the plant cells observed. Therefore, the polynucleotides and polypeptides can be employed to change expression levels of a genes, polynucleotides, and/or proteins of plants.

Producing Polypeptides

The polynucleotides of the invention include sequences that encode transcription factors and transcription factor homolog polypeptides and sequences complementary thereto, as well as unique fragments of coding sequence, or sequence complementary thereto. Such polynucleotides can be, e.g., DNA or RNA, e.g., mRNA, cRNA, synthetic RNA, genomic DNA, cDNA synthetic DNA, oligonucleotides, etc. The polynucleotides are either double-stranded or single-stranded, and include either, or both sense (i.e., coding) sequences and antisense (i.e., non-coding, complementary) sequences. The polynucleotides include the coding sequence of a transcription factor, or transcription factor homolog polypeptide, in isolation, in combination with additional coding sequences (e.g., a purification tag, a localization signal, as a fusion-protein, as a pre-protein, or the like), in combination with non-coding sequences (e.g., introns or inteins, regulatory elements such as promoters, enhancers, terminators, and the like), and/or in a vector or host environment in which the polynucleotide encoding a transcription factor or transcription factor homolog polypeptide is an endogenous or exogenous gene.

A variety of methods exist for producing the polynucleotides of the invention. Procedures for identifying and isolating DNA clones are well known to those of skill in the art and are described in, e.g., Berger and Kimmel, Guide to Molecular Cloning Techniques, Methods in Enzymology, vol. 152 Academic Press, Inc., San Diego, CA ("Berger"); Sambrook et al. (1989) Molecular Cloning - A Laboratory Manual (2nd Edition), Vol. 1-3, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, and Current Protocols in Molecular Biology, Ausubel et al. editors, Current Protocols, Greene Publishing Associates, Inc. and John Wiley & Sons, Inc. (supplemented through 2000) ("Ausubel").

Alternatively, polynucleotides of the invention, can be produced by a variety of in vitro amplification methods adapted to the present invention by appropriate selection of specific or degenerate primers. Examples of protocols sufficient to direct persons of skill through in vitro amplification methods, including the polymerase chain reaction (PCR) the ligase chain reaction (LCR), Q β -replicase amplification and other RNA polymerase mediated techniques (e.g., NASBA), e.g., for the production of

the homologous nucleic acids of the invention are found in Berger (supra), Sambrook (supra), and Ausubel (supra), as well as Mullis et al. (1987) PCR Protocols A Guide to Methods and Applications (Innis et al. eds) Academic Press Inc. San Diego, CA (1990) (Innis). Improved methods for cloning in vitro amplified nucleic acids are described in Wallace et al. U.S. Pat. No. 5,426,039. Improved methods for amplifying large nucleic acids by PCR are summarized in Cheng et al. (1994) *Nature* 369: 684-685 and the references cited therein, in which PCR amplicons of up to 40kb are generated. One of skill will appreciate that essentially any RNA can be converted into a double stranded DNA suitable for restriction digestion, PCR expansion and sequencing using reverse transcriptase and a polymerase. See, e.g., Ausubel, Sambrook and Berger, all *supra*.

Alternatively, polynucleotides and oligonucleotides of the invention can be assembled from fragments produced by solid-phase synthesis methods. Typically, fragments of up to approximately 100 bases are individually synthesized and then enzymatically or chemically ligated to produce a desired sequence, e.g., a polynucleotide encoding all or part of a transcription factor. For example, chemical synthesis using the phosphoramidite method is described, e.g., by Beaucage et al. (1981) *Tetrahedron Letters* 22: 1859-1869; and Matthes et al. (1984) *EMBO J.* 3: 801-805. According to such methods, oligonucleotides are synthesized, purified, annealed to their complementary strand, ligated and then optionally cloned into suitable vectors. And if so desired, the polynucleotides and polypeptides of the invention can be custom ordered from any of a number of commercial suppliers.

Homologous Sequences

Sequences homologous to those provided in the Sequence Listing derived from *Arabidopsis thaliana* or from other plants of choice, are also an aspect of the invention. Homologous sequences can be derived from any plant including monocots and dicots and in particular agriculturally important plant species, including but not limited to, crops such as soybean, wheat, corn (maize), potato, cotton, rice, rape, oilseed rape (including canola), sunflower, alfalfa, clover, sugarcane, and turf; or fruits and vegetables, such as banana, blackberry, blueberry, strawberry, and raspberry, cantaloupe, carrot, cauliflower, coffee, cucumber, eggplant, grapes, honeydew, lettuce, mango, melon, onion, papaya, peas, peppers, pineapple, pumpkin, spinach, squash, sweet corn, tobacco, tomato, tomatillo, watermelon, rosaceous fruits (such as apple, peach, pear, cherry and plum) and vegetable brassicas (such as broccoli, cabbage, cauliflower, Brussels sprouts, and kohlrabi). Other crops, including fruits and vegetables, whose phenotype can be changed and which comprise homologous sequences include barley; rye; millet; sorghum; currant; avocado; citrus fruits such as oranges, lemons, grapefruit and tangerines, artichoke, cherries; nuts such as the walnut and peanut; endive; leek; roots such as arrowroot, beet, cassava, turnip, radish, yam, and sweet

potato; and beans. The homologous sequences may also be derived from woody species, such as pine, poplar and eucalyptus, or mint or other labiates. In addition, homologous sequences may be derived from plants that are evolutionarily-related to crop plants, but which may not have yet been used as crop plants. Examples include deadly nightshade (*Atropa belladonna*), related to tomato; jimson weed (*Datura* 5 *strammium*), related to peyote; and teosinte (*Zea* species), related to corn (maize).

Orthologs and Paralogs

Homologous sequences as described above can comprise orthologous or paralogous sequences. Several different methods are known by those of skill in the art for identifying and defining these 10 functionally homologous sequences. Three general methods for defining orthologs and paralogs are described; an ortholog or paralog, including equivalogs, may be identified by one or more of the methods described below.

Orthologs and paralogs are evolutionarily related genes that have similar sequence and similar functions. Orthologs are structurally related genes in different species that are derived by a speciation 15 event. Paralogs are structurally related genes within a single species that are derived by a duplication event.

Within a single plant species, gene duplication may cause two copies of a particular gene, giving rise to two or more genes with similar sequence and often similar function known as paralogs. A paralog is therefore a similar gene formed by duplication within the same species. Paralogs typically cluster together 20 or in the same clade (a group of similar genes) when a gene family phylogeny is analyzed using programs such as CLUSTAL (Thompson et al. (1994) *Nucleic Acids Res.* 22: 4673-4680; Higgins et al. (1996) *Methods Enzymol.* 266: 383-402). Groups of similar genes can also be identified with pair-wise BLAST analysis (Feng and Doolittle (1987) *J. Mol. Evol.* 25: 351-360). For example, a clade of very similar MADS domain transcription factors from *Arabidopsis* all share a common function in flowering time 25 (Ratcliffe et al. (2001) *Plant Physiol.* 126: 122-132), and a group of very similar AP2 domain transcription factors from *Arabidopsis* are involved in tolerance of plants to freezing (Gilmour et al. (1998) *Plant J.* 16: 433-442). Analysis of groups of similar genes with similar function that fall within one clade can yield sub-sequences that are particular to the clade. These sub-sequences, known as consensus sequences, can not only be used to define the sequences within each clade, but define the 30 functions of these genes; genes within a clade may contain paralogous sequences, or orthologous sequences that share the same function (see also, for example, Mount (2001), in Bioinformatics: Sequence and Genome Analysis, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, page 543.)

Speciation, the production of new species from a parental species, can also give rise to two or

more genes with similar sequence and similar function. These genes, termed orthologs, often have an identical function within their host plants and are often interchangeable between species without losing function. Because plants have common ancestors, many genes in any plant species will have a corresponding orthologous gene in another plant species. Once a phylogenetic tree for a gene family of one species has been constructed using a program such as CLUSTAL (Thompson et al. (1994) *Nucleic Acids Res.* 22: 4673-4680; Higgins et al. (1996) *supra*) potential orthologous sequences can be placed into the phylogenetic tree and their relationship to genes from the species of interest can be determined. Orthologous sequences can also be identified by a reciprocal BLAST strategy. Once an orthologous sequence has been identified, the function of the ortholog can be deduced from the identified function of the reference sequence.

Transcription factor gene sequences are conserved across diverse eukaryotic species lines (Goodrich et al. (1993) *Cell* 75: 519-530; Lin et al. (1991) *Nature* 353: 569-571; Sadowski et al. (1988) *Nature* 335: 563-564). Plants are no exception to this observation; diverse plant species possess transcription factors that have similar sequences and functions.

Orthologous genes from different organisms have highly conserved functions, and very often essentially identical functions (Lee et al. (2002) *Genome Res.* 12: 493-502; Remm et al. (2001) *J. Mol. Biol.* 314: 1041-1052). Paralogous genes, which have diverged through gene duplication, may retain similar functions of the encoded proteins. In such cases, paralogs can be used interchangeably with respect to certain embodiments of the instant invention (for example, transgenic expression of a coding sequence). An example of such highly related paralogs is the CBF family, with three well-defined members in *Arabidopsis* and at least one ortholog in *Brassica napus* (SEQ ID NOs: 69, 71, 73, or 75, respectively), all of which control pathways involved in both freezing and drought stress (Gilmour et al. (1998) *Plant J.* 16: 433-442; Jaglo et al. (1998) *Plant Physiol.* 127: 910-917).

The following references represent a small sampling of the many studies that demonstrate that conserved transcription factor genes from diverse species are likely to function similarly (i.e., regulate similar target sequences and control the same traits), and that transcription factors may be transformed into diverse species to confer or improve traits.

(1) The *Arabidopsis* NPR1 gene regulates systemic acquired resistance (SAR); over-expression of NPR1 leads to enhanced resistance in *Arabidopsis*. When either *Arabidopsis* NPR1 or the rice NPR1 ortholog was overexpressed in rice (which, as a monocot, is diverse from *Arabidopsis*), challenge with the rice bacterial blight pathogen *Xanthomonas oryzae* pv. *Oryzae*, the transgenic plants displayed enhanced resistance (Chern et al. (2001) *Plant J.* 27: 101-113). NPR1 acts

through activation of expression of transcription factor genes, such as TGA2 (Fan and Dong (2002) *Plant Cell* 14: 1377-1389).

- (2) E2F genes are involved in transcription of plant genes for proliferating cell nuclear antigen (PCNA). Plant E2Fs share a high degree of similarity in amino acid sequence between monocots and dicots, and are even similar to the conserved domains of the animal E2Fs. Such conservation indicates a functional similarity between plant and animal E2Fs. E2F transcription factors that regulate meristem development act through common cis-elements, and regulate related (PCNA) genes. (Kosugi and Ohashi, (2002) *Plant J.* 29: 45-59.)
- (3) The ABI5 gene (abscisic acid (ABA) insensitive 5) encodes a basic leucine zipper factor required for ABA response in the seed and vegetative tissues. Co-transformation experiments with ABI5 cDNA constructs in rice protoplasts resulted in specific transactivation of the ABA-inducible wheat, *Arabidopsis*, bean, and barley promoters. These results demonstrate that sequentially similar ABI5 transcription factors are key targets of a conserved ABA signaling pathway in diverse plants. (Gampala et al. (2001) *J. Biol. Chem.* 277: 1689-1694.)
- (4) Sequences of three *Arabidopsis* GAMYB-like genes were obtained on the basis of sequence similarity to GAMYB genes from barley, rice, and *L. temulentum*. These three *Arabidopsis* genes were determined to encode transcription factors (AtMYB33, AtMYB65, and AtMYB101) and could substitute for a barley GAMYB and control alpha-amylase expression. (Gocal et al. (2001) *Plant Physiol.* 127: 1682-1693.)
- (5) The floral control gene LEAFY from *Arabidopsis* can dramatically accelerate flowering in numerous dicotyledonous plants. Constitutive expression of *Arabidopsis* LEAFY also caused early flowering in transgenic rice (a monocot), with a heading date that was 26-34 days earlier than that of wild-type plants. These observations indicate that floral regulatory genes from *Arabidopsis* are useful tools for heading date improvement in cereal crops. (He et al. (2000) *Transgenic Res.* 9: 223-227.)
- (6) Bioactive gibberellins (GAs) are essential endogenous regulators of plant growth. GA signaling tends to be conserved across the plant kingdom. GA signaling is mediated via GAI, a nuclear member of the GRAS family of plant transcription factors. *Arabidopsis* GAI has been shown to function in rice to inhibit gibberellin response pathways. (Fu et al. (2001) *Plant Cell* 13: 1791-1802.)
- (7) The *Arabidopsis* gene SUPERMAN (SUP), encodes a putative transcription factor that maintains the boundary between stamens and carpels. By over-expressing *Arabidopsis* SUP in rice, the effect of the gene's presence on whorl boundaries was shown to be conserved. This demonstrated

that SUP is a conserved regulator of floral whorl boundaries and affects cell proliferation. (Nandi et al. (2000) *Curr. Biol.* 10: 215-218.)

(8) Maize, petunia and *Arabidopsis* myb transcription factors that regulate flavonoid biosynthesis are very genetically similar and affect the same trait in their native species, therefore sequence and function of these myb transcription factors correlate with each other in these diverse species (Borevitz et al. (2000) *Plant Cell* 12: 2383-2394).

(9) Wheat reduced height-1 (Rht-B1/Rht-D1) and maize dwarf-8 (d8) genes are orthologs of the *Arabidopsis* gibberellin insensitive (GAI) gene. Both of these genes have been used to produce dwarf grain varieties that have improved grain yield. These genes encode proteins that resemble nuclear transcription factors and contain an SH2-like domain, indicating that phosphotyrosine may participate in gibberellin signaling. Transgenic rice plants containing a mutant GAI allele from *Arabidopsis* have been shown to produce reduced responses to gibberellin and are dwarfed, indicating that mutant GAI orthologs could be used to increase yield in a wide range of crop species. (Peng et al. (1999) *Nature* 400: 256-261.)

Transcription factors that are homologous to the listed AT-hook transcription factors will typically share at least about 78% and 62% amino acid sequence identity in their AT-hook and second conserved domains, respectively. More closely related transcription factors can share at least about 89% or about 100% identity in their AT-hook domains, and at least about 63%, or at least about 65%, or at least about 67%, or at least about 68%, or at least about 69%, or at least about 71%, or at least about 100% identity with the second conserved domain of G1073, as seen by the examples shown to have function in Table 1.. At the nucleotide level, the sequences of the invention will typically share at least about 40% nucleotide sequence identity, preferably at least about 50%, about 60%, about 70% or about 80% sequence identity, and more preferably about 85%, about 90%, about 95% or about 97% or more sequence identity to one or more of the listed full-length sequences, or to a listed sequence but excluding or outside a known consensus sequence or consensus DNA-binding site, or outside one or all conserved domain. The degeneracy of the genetic code enables major variations in the nucleotide sequence of a polynucleotide while maintaining the amino acid sequence of the encoded protein. Conserved domains within the AT-hook transcription factor family may exhibit a higher degree of sequence homology, such as at least 62% amino acid sequence identity including conservative substitutions, and preferably at least 65% sequence identity, and more preferably at least 69%, or at least about 71%, or at least about 78%, or at least about 89%, or at least about 90%, or at least about 95%, or at least about 98% sequence identity. Transcription factors that are homologous to the listed sequences should share at least 50%, or at least about 60%, or at

least about 75%, or at least about 80%, or at least about 90%, or at least about 95% amino acid sequence identity over the entire length of the polypeptide or the homolog.

Percent identity can be determined electronically, e.g., by using the MEGALIGN program (DNASTAR, Inc. Madison, Wis.). The MEGALIGN program can create alignments between two or more sequences according to different methods, for example, the clustal method. (See, for example, Higgins and Sharp (1988) *Gene* 73: 237-244.) The clustal algorithm groups sequences into clusters by examining the distances between all pairs. The clusters are aligned pairwise and then in groups. Other alignment algorithms or programs may be used, including FASTA, BLAST, or ENTREZ, FASTA and BLAST, and which may be used to calculate percent similarity. These are available as a part of the GCG sequence analysis package (University of Wisconsin, Madison, Wis.), and can be used with or without default settings. ENTREZ is available through the National Center for Biotechnology Information. In one embodiment, the percent identity of two sequences can be determined by the GCG program with a gap weight of 1, e.g., each amino acid gap is weighted as if it were a single amino acid or nucleotide mismatch between the two sequences (see USPN 6,262,333).

Other techniques for alignment are described in *Methods in Enzymology*, vol. 266, Computer Methods for Macromolecular Sequence Analysis (1996), ed. Doolittle, Academic Press, Inc., San Diego, Calif., USA. Preferably, an alignment program that permits gaps in the sequence is utilized to align the sequences. The Smith-Waterman is one type of algorithm that permits gaps in sequence alignments (see Shpaer (1997) *Methods Mol. Biol.* 70: 173-187). Also, the GAP program using the Needleman and Wunsch alignment method can be utilized to align sequences. An alternative search strategy uses MPSRCH software, which runs on a MASPAR computer. MPSRCH uses a Smith-Waterman algorithm to score sequences on a massively parallel computer. This approach improves ability to pick up distantly related matches, and is especially tolerant of small gaps and nucleotide sequence errors. Nucleic acid-encoded amino acid sequences can be used to search both protein and DNA databases.

The percentage similarity between two polypeptide sequences, e.g., sequence A and sequence B, is calculated by dividing the length of sequence A, minus the number of gap residues in sequence A, minus the number of gap residues in sequence B, into the sum of the residue matches between sequence A and sequence B, times one hundred. Gaps of low or of no similarity between the two amino acid sequences are not included in determining percentage similarity. Percent identity between polynucleotide sequences can also be counted or calculated by other methods known in the art, e.g., the Jotun Hein method. (See, for example, Hein (1990) *Methods Enzymol.* 183: 626-645.) Identity between sequences can also be determined by other methods known in the art, e.g., by varying hybridization conditions (see US Patent Application No. 20010010913).

Thus, the invention provides methods for identifying a sequence similar or paralogous or orthologous or homologous to one or more polynucleotides as noted herein, or one or more target polypeptides encoded by the polynucleotides, or otherwise noted herein and may include linking or associating a given plant phenotype or gene function with a sequence. In the methods, a sequence database is provided (locally or across an internet or intranet) and a query is made against the sequence database using the relevant sequences herein and associated plant phenotypes or gene functions.

In addition, one or more polynucleotide sequences or one or more polypeptides encoded by the polynucleotide sequences may be used to search against a BLOCKS (Bairoch et al. (1997) *Nucleic Acids Res.* 25: 217-221), PFAM, and other databases which contain previously identified and annotated motifs, sequences and gene functions. Methods that search for primary sequence patterns with secondary structure gap penalties (Smith et al. (1992) *Protein Engineering* 5: 35-51) as well as algorithms such as Basic Local Alignment Search Tool (BLAST; Altschul (1993) *J. Mol. Evol.* 36: 290-300; Altschul et al. (1990) *supra*), BLOCKS (Henikoff and Henikoff (1991) *Nucleic Acids Res.* 19: 6565-6572), Hidden Markov Models (HMM; Eddy (1996) *Curr. Opin. Str. Biol.* 6: 361-365; Sonnhammer et al. (1997) *Proteins* 28: 405-420), and the like, can be used to manipulate and analyze polynucleotide and polypeptide sequences encoded by polynucleotides. These databases, algorithms and other methods are well known in the art and are described in Ausubel et al. (1997; Short Protocols in Molecular Biology, John Wiley & Sons, New York, NY, unit 7.7) and in Meyers (1995; Molecular Biology and Biotechnology, Wiley VCH, New York, NY, p 856-853).

A further method for identifying or confirming that specific homologous sequences control the same function is by comparison of the transcript profile(s) obtained upon overexpression or knockout of two or more related transcription factors. Since transcript profiles are diagnostic for specific cellular states, one skilled in the art will appreciate that genes that have a highly similar transcript profile (e.g., with greater than 50% regulated transcripts in common, more preferably with greater than 70% regulated transcripts in common, most preferably with greater than 90% regulated transcripts in common) will have highly similar functions. Fowler et al. (2002, *Plant Cell*, 14: 1675-1679) have shown that three paralogous AP2 family genes (CBF1, CBF2 and CBF3), each of which is induced upon cold treatment, and each of which can condition improved freezing tolerance, have highly similar transcript profiles. Once a transcription factor has been shown to provide a specific function, its transcript profile becomes a diagnostic tool to determine whether putative paralogs or orthologs have the same function.

Furthermore, methods using manual alignment of sequences similar or homologous to one or more polynucleotide sequences or one or more polypeptides encoded by the polynucleotide sequences may be used to identify regions of similarity and AT-hook domains. Such manual methods are well-known of

those of skill in the art and can include, for example, comparisons of tertiary structure between a polypeptide sequence encoded by a polynucleotide which comprises a known function with a polypeptide sequence encoded by a polynucleotide sequence which has a function not yet determined. Such examples of tertiary structure may comprise predicted alpha helices, beta-sheets, amphipathic helices, leucine zipper motifs, zinc finger motifs, proline-rich regions, cysteine repeat motifs, and the like.

Orthologs and paralogs of presently disclosed transcription factors may be cloned using compositions provided by the present invention according to methods well known in the art. cDNAs can be cloned using mRNA from a plant cell or tissue that expresses one of the present transcription factors. Appropriate mRNA sources may be identified by interrogating Northern blots with probes designed from the present transcription factor sequences, after which a library is prepared from the mRNA obtained from a positive cell or tissue. Transcription factor-encoding cDNA is then isolated using, for example, PCR, using primers designed from a presently disclosed transcription factor gene sequence, or by probing with a partial or complete cDNA or with one or more sets of degenerate probes based on the disclosed sequences. The cDNA library may be used to transform plant cells. Expression of the cDNAs of interest is detected using, for example, methods disclosed herein such as microarrays, Northern blots, quantitative PCR, or any other technique for monitoring changes in expression. Genomic clones may be isolated using similar techniques to those.

Examples of orthologs of the *Arabidopsis* polypeptide sequences SEQ ID NOs: 2, 4, 6, and 8 include SEQ ID NOs: 10, 12, 14, 16, 18, and other functionally similar orthologs listed in the Sequence Listing. In addition to the sequences in the Sequence Listing, the invention encompasses isolated nucleotide sequences that are sequentially and structurally similar to G1073, G1067, G2153, G2156, G3399, G3407, G3456, G3459 and G3460 (SEQ ID NO: 1, 3, 5, 7, 9, 11, 13, 15, 17) and function in a plant by increasing biomass and regulating abiotic stress tolerance. These polypeptide sequences show sequence similarity to G1073, as shown by their respective identities to G1073 and the conserved domains of G1073, in Table 1.

Since all of these polynucleotide sequences are phylogenetically related and similar in sequence (the phylogenetic tree shown in Figure 3 includes many of these sequences), and have been shown to increase a plant's biomass, one skilled in the art would predict that other similar, phylogenetically related sequences would also increase a plant's biomass. Since a number of these structurally related sequences have also been shown to increase abiotic stress tolerance, one skilled in the art would conclude that phylogenetically related equivalents of these sequences would function in a similar capacity.

Identifying Polynucleotides or Nucleic Acids by Hybridization

Polynucleotides homologous to the sequences illustrated in the Sequence Listing and tables can be identified, e.g., by hybridization to each other under stringent or under highly stringent conditions. Single stranded polynucleotides hybridize when they associate based on a variety of well characterized physical-chemical forces, such as hydrogen bonding, solvent exclusion, base stacking and the like. The stringency of a hybridization reflects the degree of sequence identity of the nucleic acids involved, such that the higher the stringency, the more similar are the two polynucleotide strands. Stringency is influenced by a variety of factors, including temperature, salt concentration and composition, organic and non-organic additives, solvents, etc. present in both the hybridization and wash solutions and incubations (and number thereof), as described in more detail in the references cited above.

Encompassed by the invention are polynucleotide sequences that are capable of hybridizing to the claimed polynucleotide sequences, including any of the transcription factor polynucleotides within the Sequence Listing, and fragments thereof under various conditions of stringency (See, for example, Wahl and Berger (1987) *Methods Enzymol.* 152: 399-407; and Kimmel (1987) *Methods Enzymol.* 152: 507-511). In addition to the nucleotide sequences listed in the Sequence Listing, full length cDNA, orthologs, and paralogs of the present nucleotide sequences may be identified and isolated using well-known methods. The cDNA libraries, orthologs, and paralogs of the present nucleotide sequences may be screened using hybridization methods to determine their utility as hybridization target or amplification probes.

With regard to hybridization, conditions that are highly stringent, and means for achieving them, are well known in the art. See, for example, Sambrook et al. (1989) "*Molecular Cloning: A Laboratory Manual*" (2nd ed., Cold Spring Harbor Laboratory); Berger and Kimmel, eds., (1987) "Guide to Molecular Cloning Techniques", In *Methods in Enzymology*:152: 467-469; and Anderson and Young (1985) "Quantitative Filter Hybridisation." In: Hames and Higgins, ed., Nucleic Acid Hybridisation, A Practical Approach. Oxford, IRL Press, 73-111.

Stability of DNA duplexes is affected by such factors as base composition, length, and degree of base pair mismatch. Hybridization conditions may be adjusted to allow DNAs of different sequence relatedness to hybridize. The melting temperature (T_m) is defined as the temperature when 50% of the duplex molecules have dissociated into their constituent single strands. The melting temperature of a perfectly matched duplex, where the hybridization buffer contains formamide as a denaturing agent, may be estimated by the following equations:

(I) DNA-DNA:

$$T_m(^{\circ}\text{C})=81.5+16.6(\log [\text{Na}^{+}])+0.41(\% \text{ G+C})-0.62(\% \text{ formamide})-500/L$$

(II) DNA-RNA:

$$5 \quad T_m(^{\circ}\text{C})=79.8+18.5(\log [\text{Na}^{+}])+0.58(\% \text{ G+C})+0.12(\% \text{ G+C})^2-0.5(\% \text{ formamide})-820/L$$

(III) RNA-RNA:

$$T_m(^{\circ}\text{C})=79.8+18.5(\log [\text{Na}^{+}])+0.58(\% \text{ G+C})+0.12(\% \text{ G+C})^2-0.35(\% \text{ formamide})-820/L$$

10 where L is the length of the duplex formed, $[\text{Na}^{+}]$ is the molar concentration of the sodium ion in the hybridization or washing solution, and % G+C is the percentage of (guanine+cytosine) bases in the hybrid. For imperfectly matched hybrids, approximately 1°C is required to reduce the melting temperature for each 1% mismatch.

Hybridization experiments are generally conducted in a buffer of pH between 6.8 to 7.4, although
 15 the rate of hybridization is nearly independent of pH at ionic strengths likely to be used in the hybridization buffer (Anderson et al. (1985) *supra*). In addition, one or more of the following may be used to reduce non-specific hybridization: sonicated salmon sperm DNA or another non-complementary DNA, bovine serum albumin, sodium pyrophosphate, sodium dodecylsulfate (SDS), polyvinylpyrrolidone, ficoll and Denhardt's solution. Dextran sulfate and polyethylene glycol 6000 act to exclude
 20 DNA from solution, thus raising the effective probe DNA concentration and the hybridization signal within a given unit of time. In some instances, conditions of even greater stringency may be desirable or required to reduce non-specific and/or background hybridization. These conditions may be created with the use of higher temperature, lower ionic strength and higher concentration of a denaturing agent such as formamide.

25 Stringency conditions can be adjusted to screen for moderately similar fragments such as homologous sequences from distantly related organisms, or to highly similar fragments such as genes that duplicate functional enzymes from closely related organisms. The stringency can be adjusted either during the hybridization step or in the post-hybridization washes. Salt concentration, formamide concentration, hybridization temperature and probe lengths are variables that can be used to alter stringency (as described
 30 by the formula above). As a general guidelines high stringency is typically performed at $T_m-5^{\circ}\text{C}$ to $T_m-20^{\circ}\text{C}$, moderate stringency at $T_m-20^{\circ}\text{C}$ to $T_m-35^{\circ}\text{C}$ and low stringency at $T_m-35^{\circ}\text{C}$ to $T_m-50^{\circ}\text{C}$ for duplex >150 base pairs. Hybridization may be performed at low to moderate stringency ($25-50^{\circ}\text{C}$ below T_m), followed by post-hybridization washes at increasing stringencies. Maximum rates of hybridization in

solution are determined empirically to occur at $T_m-25^\circ\text{C}$ for DNA-DNA duplex and $T_m-15^\circ\text{C}$ for RNA-DNA duplex. Optionally, the degree of dissociation may be assessed after each wash step to determine the need for subsequent, higher stringency wash steps.

High stringency conditions may be used to select for nucleic acid sequences with high degrees of identity to the disclosed sequences. An example of stringent hybridization conditions obtained in a filter-based method such as a Southern or northern blot for hybridization of complementary nucleic acids that have more than 100 complementary residues is about 5°C to 20°C lower than the thermal melting point (T_m) for the specific sequence at a defined ionic strength and pH. Conditions used for hybridization may include about 0.02 M to about 0.15 M sodium chloride, about 0.5% to about 5% casein, about 0.02% SDS or about 0.1% N-laurylsarcosine, about 0.001 M to about 0.03 M sodium citrate, at hybridization temperatures between about 50°C and about 70°C . More preferably, high stringency conditions are about 0.02 M sodium chloride, about 0.5% casein, about 0.02% SDS, about 0.001 M sodium citrate, at a temperature of about 50°C . Nucleic acid molecules that hybridize under stringent conditions will typically hybridize to a probe based on either the entire DNA molecule or selected portions, e.g., to a unique subsequence, of the DNA.

Stringent salt concentration will ordinarily be less than about 750 mM NaCl and 75 mM trisodium citrate. Increasingly stringent conditions may be obtained with less than about 500 mM NaCl and 50 mM trisodium citrate, to even greater stringency with less than about 250 mM NaCl and 25 mM trisodium citrate. Low stringency hybridization can be obtained in the absence of organic solvent, e.g., formamide, whereas high stringency hybridization may be obtained in the presence of at least about 35% formamide, and more preferably at least about 50% formamide. Stringent temperature conditions will ordinarily include temperatures of at least about 30°C , more preferably of at least about 37°C , and most preferably of at least about 42°C with formamide present. Varying additional parameters, such as hybridization time, the concentration of detergent, e.g., sodium dodecyl sulfate (SDS) and ionic strength, are well known to those skilled in the art. Various levels of stringency are accomplished by combining these various conditions as needed.

The washing steps that follow hybridization may also vary in stringency; the post-hybridization wash steps primarily determine hybridization specificity, with the most critical factors being temperature and the ionic strength of the final wash solution. Wash stringency can be increased by decreasing salt concentration or by increasing temperature. Stringent salt concentration for the wash steps will preferably be less than about 30 mM NaCl and 3 mM trisodium citrate, and most preferably less than about 15 mM NaCl and 1.5 mM trisodium citrate.

Thus, hybridization and wash conditions that may be used to bind and remove polynucleotides

with less than the desired homology to the nucleic acid sequences or their complements that encode the present transcription factors include, for example:

6X SSC at 65° C;

50% formamide, 4X SSC at 42° C; or

5 0.5X SSC, 0.1% SDS at 65° C;

with, for example, two wash steps of 10 - 30 minutes each. . Useful variations on these conditions will be readily apparent to those skilled in the art.

A person of skill in the art would not expect substantial variation among polynucleotide species encompassed within the scope of the present invention because the highly stringent conditions set forth in
10 the above formulae yield structurally similar polynucleotides.

If desired, one may employ wash steps of even greater stringency, including about 0.2x SSC, 0.1% SDS at 65° C and washing twice, each wash step being about 30 min, or about 0.1 x SSC, 0.1% SDS at 65° C and washing twice for 30 min. The temperature for the wash solutions will ordinarily be at least about 25° C, and for greater stringency at least about 42° C. Hybridization stringency may be increased
15 further by using the same conditions as in the hybridization steps, with the wash temperature raised about 3° C to about 5° C, and stringency may be increased even further by using the same conditions except the wash temperature is raised about 6° C to about 9° C. For identification of less closely related homologs, wash steps may be performed at a lower temperature, e.g., 50° C.

An example of a low stringency wash step employs a solution and conditions of at least 25° C in
20 30 mM NaCl, 3 mM trisodium citrate, and 0.1% SDS over 30 min. Greater stringency may be obtained at 42° C in 15 mM NaCl, with 1.5 mM trisodium citrate, and 0.1% SDS over 30 min. Even higher stringency wash conditions are obtained at 65° C -68° C in a solution of 15 mM NaCl, 1.5 mM trisodium citrate, and 0.1% SDS. Wash procedures will generally employ at least two final wash steps. Additional variations on these conditions will be readily apparent to those skilled in the art (see, for example, US
25 Patent Application No. 20010010913).

Stringency conditions can be selected such that an oligonucleotide that is perfectly complementary to the coding oligonucleotide hybridizes to the coding oligonucleotide with at least about a 5-10x higher signal to noise ratio than the ratio for hybridization of the perfectly complementary oligonucleotide to a nucleic acid encoding a transcription factor known as of the filing date of the application. It may be
30 desirable to select conditions for a particular assay such that a higher signal to noise ratio, that is, about 15x or more, is obtained. Accordingly, a subject nucleic acid will hybridize to a unique coding oligonucleotide with at least a 2x or greater signal to noise ratio as compared to hybridization of the coding oligonucleotide to a nucleic acid encoding known polypeptide. The particular signal will depend

on the label used in the relevant assay, e.g., a fluorescent label, a colorimetric label, a radioactive label, or the like. Labeled hybridization or PCR probes for detecting related polynucleotide sequences may be produced by oligolabeling, nick translation, end-labeling, or PCR amplification using a labeled nucleotide.

Encompassed by the invention are polynucleotide sequences that are capable of hybridizing to the claimed polynucleotide sequences, including any of the transcription factor polynucleotides within the Sequence Listing, and fragments thereof under various conditions of stringency (See, for example, Wahl and Berger (1987) *Methods Enzymol.* 152: 399-407; and Kimmel (1987) *Methods Enzymol.* 152: 507-511). In addition to the nucleotide sequences in the Sequence Listing, full length cDNA, orthologs, and paralogs of the present nucleotide sequences may be identified and isolated using well-known methods. The cDNA libraries, orthologs, and paralogs of the present nucleotide sequences may be screened using hybridization methods to determine their utility as hybridization target or amplification probes.

Identifying Polynucleotides or Nucleic Acids with Expression Libraries

In addition to hybridization methods, transcription factor homolog polypeptides can be obtained by screening an expression library using antibodies specific for one or more transcription factors. With the provision herein of the disclosed transcription factor, and transcription factor homolog nucleic acid sequences, the encoded polypeptide(s) can be expressed and purified in a heterologous expression system (for example, *E. coli*) and used to raise antibodies (monoclonal or polyclonal) specific for the polypeptide(s) in question. Antibodies can also be raised against synthetic peptides derived from transcription factor, or transcription factor homolog, amino acid sequences. Methods of raising antibodies are well known in the art and are described in Harlow and Lane (1988), Antibodies: A Laboratory Manual, Cold Spring Harbor Laboratory, New York. Such antibodies can then be used to screen an expression library produced from the plant from which it is desired to clone additional transcription factor homologs, using the methods described above. The selected cDNAs can be confirmed by sequencing and enzymatic activity.

Sequence Variations

It will readily be appreciated by those of skill in the art, that any of a variety of polynucleotide sequences are capable of encoding the transcription factors and transcription factor homolog polypeptides of the invention. Due to the degeneracy of the genetic code, many different polynucleotides can encode identical and/or substantially similar polypeptides in addition to those sequences illustrated in the Sequence Listing. Nucleic acids having a sequence that differs from the sequences shown in the Sequence Listing, or complementary sequences, that encode functionally equivalent peptides (i.e., peptides having

some degree of equivalent or similar biological activity) but differ in sequence from the sequence shown in the Sequence Listing due to degeneracy in the genetic code, are also within the scope of the invention.

Altered polynucleotide sequences encoding polypeptides include those sequences with deletions, insertions, or substitutions of different nucleotides, resulting in a polynucleotide encoding a polypeptide with at least one functional characteristic of the instant polypeptides. Included within this definition are polymorphisms which may or may not be readily detectable using a particular oligonucleotide probe of the polynucleotide encoding the instant polypeptides, and improper or unexpected hybridization to allelic variants, with a locus other than the normal chromosomal locus for the polynucleotide sequence encoding the instant polypeptides.

Allelic variant refers to any of two or more alternative forms of a gene occupying the same chromosomal locus. Allelic variation arises naturally through mutation, and may result in phenotypic polymorphism within populations. Gene mutations can be silent (i.e., no change in the encoded polypeptide) or may encode polypeptides having altered amino acid sequence. The term allelic variant is also used herein to denote a protein encoded by an allelic variant of a gene. Splice variant refers to alternative forms of RNA transcribed from a gene. Splice variation arises naturally through use of alternative splicing sites within a transcribed RNA molecule, or less commonly between separately transcribed RNA molecules, and may result in several mRNAs transcribed from the same gene. Splice variants may encode polypeptides having altered amino acid sequence. The term splice variant is also used herein to denote a protein encoded by a splice variant of an mRNA transcribed from a gene.

Those skilled in the art would recognize that, for example, G1073, SEQ ID NO: 2, represents a single transcription factor; allelic variation and alternative splicing may be expected to occur. Allelic variants of SEQ ID NO: 1 can be cloned by probing cDNA or genomic libraries from different individual organisms according to standard procedures. Allelic variants of the DNA sequence shown in SEQ ID NO: 1, including those containing silent mutations and those in which mutations result in amino acid sequence changes, are within the scope of the present invention, as are proteins which are allelic variants of SEQ ID NO: 2. cDNAs generated from alternatively spliced mRNAs, which retain the properties of the transcription factor are included within the scope of the present invention, as are polypeptides encoded by such cDNAs and mRNAs. Allelic variants and splice variants of these sequences can be cloned by probing cDNA or genomic libraries from different individual organisms or tissues according to standard procedures known in the art (see USPN 6,388,064).

Thus, in addition to the sequences set forth in the Sequence Listing, the invention also encompasses related nucleic acid molecules that include allelic or splice variants, and sequences that are complementary. Related nucleic acid molecules also include nucleotide sequences encoding a polypeptide

comprising or consisting essentially of a substitution, modification, addition and/or deletion of one or more amino acid residues. Such related polypeptides may comprise, for example, additions and/or deletions of one or more N-linked or O-linked glycosylation sites, or an addition and/or a deletion of one or more cysteine residues.

5 For example, Table 2 illustrates, for example, that the codons AGC, AGT, TCA, TCC, TCG, and TCT all encode the same amino acid: serine. Accordingly, at each position in the sequence where there is a codon encoding serine, any of the above trinucleotide sequences can be used without altering the encoded polypeptide.

Table 2

Amino acid			Possible Codons							
Alanine	Ala	A	GCA	GCC	GCG	GCU				
Cysteine	Cys	C	TGC	TGT						
Aspartic acid	Asp	D	GAC	GAT						
Glutamic acid	Glu	E	GAA	GAG						
Phenylalanine	Phe	F	TTC	TTT						
Glycine	Gly	G	GGA	GGC	GGG	GGT				
Histidine	His	H	CAC	CAT						
Isoleucine	Ile	I	ATA	ATC	ATT					
Lysine	Lys	K	AAA	AAG						
Leucine	Leu	L	TTA	TTG	CTA	CTC	CTG	CTT		
Methionine	Met	M	ATG							
Asparagine	Asn	N	AAC	AAT						
Proline	Pro	P	CCA	CCC	CCG	CCT				
Glutamine	Gln	Q	CAA	CAG						
Arginine	Arg	R	AGA	AGG	CGA	CGC	CGG	CGT		
Serine	Ser	S	AGC	AGT	TCA	TCC	TCG	TCT		
Threonine	Thr	T	ACA	ACC	ACG	ACT				
Valine	Val	V	GTA	GTC	GTG	GTT				
Tryptophan	Trp	W	TGG							
Tyrosine	Tyr	Y	TAC	TAT						

Sequence alterations that do not change the amino acid sequence encoded by the polynucleotide are termed "silent" variations. With the exception of the codons ATG and TGG, encoding methionine and tryptophan, respectively, any of the possible codons for the same amino acid can be substituted by a variety of techniques, e.g., site-directed mutagenesis, available in the art. Accordingly, any and all such variations of a sequence selected from the above table are a feature of the invention.

In addition to silent variations, other conservative variations that alter one, or a few amino acids in the encoded polypeptide, can be made without altering the function of the polypeptide, these conservative variants are, likewise, a feature of the invention.

For example, substitutions, deletions and insertions introduced into the sequences provided in the Sequence Listing, are also envisioned by the invention. Such sequence modifications can be engineered into a sequence by site-directed mutagenesis (Wu, editor; *Methods Enzymol.* (1993) vol. 217, Academic Press) or the other methods noted below. Amino acid substitutions are typically of single residues; insertions usually will be on the order of about from 1 to 10 amino acid residues; and deletions will range about from 1 to 30 residues. In preferred embodiments, deletions or insertions are made in adjacent pairs, e.g., a deletion of two residues or insertion of two residues. Substitutions, deletions, insertions or any combination thereof can be combined to arrive at a sequence. The mutations that are made in the polynucleotide encoding the transcription factor should not place the sequence out of reading frame and should not create complementary regions that could produce secondary mRNA structure. Preferably, the polypeptide encoded by the DNA performs the desired function.

Conservative substitutions are those in which at least one residue in the amino acid sequence has been removed and a different residue inserted in its place. Such substitutions generally are made in accordance with the Table 3 when it is desired to maintain the activity of the protein. Table 3 shows amino acids which can be substituted for an amino acid in a protein and which are typically regarded as conservative substitutions. In one embodiment, transcriptions factors listed in the Sequence Listing may have up to 10 conservative substitutions and retain their function. In another embodiment, transcriptions factors listed in the Sequence Listing may have more than 10 conservative substitutions and still retain their function.

Table 3

Residue	Conservative Substitutions
Ala	Ser
Arg	Lys
Asn	Gln; His
Asp	Glu
Gln	Asn
Cys	Ser
Glu	Asp
Gly	Pro
His	Asn; Gln
Ile	Leu, Val
Leu	Ile; Val
Lys	Arg; Gln
Met	Leu; Ile
Phe	Met; Leu; Tyr
Ser	Thr; Gly
Thr	Ser; Val
Trp	Tyr
Tyr	Trp; Phe
Val	Ile; Leu

Similar substitutions are those in which at least one residue in the amino acid sequence has been removed and a different residue inserted in its place. Such substitutions generally are made in accordance with the Table 4 when it is desired to maintain the activity of the protein. Table 4 shows amino acids which can be substituted for an amino acid in a protein and which are typically regarded as structural and functional substitutions. For example, a residue in column 1 of Table 4 may be substituted with a residue in column 2; in addition, a residue in column 2 of Table 4 may be substituted with the residue of column

10 1.

Table 4

Residue	Similar Substitutions
Ala	Ser; Thr; Gly; Val; Leu; Ile
Arg	Lys; His; Gly
Asn	Gln; His; Gly; Ser; Thr
Asp	Glu, Ser; Thr
Gln	Asn; Ala
Cys	Ser; Gly
Glu	Asp
Gly	Pro; Arg
His	Asn; Gln; Tyr; Phe; Lys; Arg
Ile	Ala; Leu; Val; Gly; Met
Leu	Ala; Ile; Val; Gly; Met
Lys	Arg; His; Gln; Gly; Pro
Met	Leu; Ile; Phe
Phe	Met; Leu; Tyr; Trp; His; Val; Ala
Ser	Thr; Gly; Asp; Ala; Val; Ile; His
Thr	Ser; Val; Ala; Gly
Trp	Tyr; Phe; His
Tyr	Trp; Phe; His
Val	Ala; Ile; Leu; Gly; Thr; Ser; Glu

Substitutions that are less conservative than those in Table 4 can be selected by picking residues that differ more significantly in their effect on maintaining (a) the structure of the polypeptide backbone in the area of the substitution, for example, as a sheet or helical conformation, (b) the charge or hydrophobicity of the molecule at the target site, or (c) the bulk of the side chain. The substitutions which in general are expected to produce the greatest changes in protein properties will be those in which (a) a hydrophilic residue, e.g., seryl or threonyl, is substituted for (or by) a hydrophobic residue, e.g., leucyl, isoleucyl, phenylalanyl, valyl or alanyl; (b) a cysteine or proline is substituted for (or by) any other residue; (c) a residue having an electropositive side chain, e.g., lysyl, arginyl, or histidyl, is substituted for (or by) an electronegative residue, e.g., glutamyl or aspartyl; or (d) a residue having a bulky side chain, e.g.,

phenylalanine, is substituted for (or by) one not having a side chain, e.g., glycine.

Further Modifying Sequences of the Invention – Mutation/Forced Evolution

In addition to generating silent or conservative substitutions as noted, above, the present invention optionally includes methods of modifying the sequences of the Sequence Listing. In the methods, nucleic acid or protein modification methods are used to alter the given sequences to produce new sequences and/or to chemically or enzymatically modify given sequences to change the properties of the nucleic acids or proteins.

Thus, in one embodiment, given nucleic acid sequences are modified, e.g., according to standard mutagenesis or artificial evolution methods to produce modified sequences. The modified sequences may be created using purified natural polynucleotides isolated from any organism or may be synthesized from purified compositions and chemicals using chemical means well known to those of skill in the art. For example, Ausubel (*supra*), provides additional details on mutagenesis methods. Artificial forced evolution methods are described, for example, by Stemmer (1994; *Nature* 370: 389-391), Stemmer (1994; *Proc. Natl. Acad. Sci.* 91: 10747-10751), and U.S. Patents 5,811,238, 5,837,500, and 6,242,568. Methods for engineering synthetic transcription factors and other polypeptides are described, for example, by Zhang et al. (2000) *J. Biol. Chem.* 275: 33850-33860, Liu et al. (2001) *J. Biol. Chem.* 276: 11323-11334, and Isalan et al. (2001) *Nature Biotechnol.* 19: 656-660. Many other mutation and evolution methods are also available and expected to be within the skill of the practitioner.

Similarly, chemical or enzymatic alteration of expressed nucleic acids and polypeptides can be performed by standard methods. For example, sequence can be modified by addition of lipids, sugars, peptides, organic or inorganic compounds, by the inclusion of modified nucleotides or amino acids, or the like. For example, protein modification techniques are illustrated in Ausubel (*supra*). Further details on chemical and enzymatic modifications can be found herein. These modification methods can be used to modify any given sequence, or to modify any sequence produced by the various mutation and artificial evolution modification methods noted herein.

Accordingly, the invention provides for modification of any given nucleic acid by mutation, evolution, chemical or enzymatic modification, or other available methods, as well as for the products produced by practicing such methods, e.g., using the sequences herein as a starting substrate for the various modification approaches.

For example, optimized coding sequence containing codons preferred by a particular prokaryotic or eukaryotic host can be used e.g., to increase the rate of translation or to produce recombinant RNA transcripts having desirable properties, such as a longer half-life, as compared with transcripts produced

using a non-optimized sequence. Translation stop codons can also be modified to reflect host preference. For example, preferred stop codons for *Saccharomyces cerevisiae* and mammals are TAA and TGA, respectively. The preferred stop codon for monocotyledonous plants is TGA, whereas insects and *E. coli* prefer to use TAA as the stop codon.

5 The polynucleotide sequences of the present invention can also be engineered in order to alter a coding sequence for a variety of reasons, including but not limited to, alterations which modify the sequence to facilitate cloning, processing and/or expression of the gene product. For example, alterations are optionally introduced using techniques which are well known in the art, e.g., site-directed mutagenesis, to insert new restriction sites, to alter glycosylation patterns, to change codon preference, to introduce
10 splice sites, etc.

Furthermore, a fragment or domain derived from any of the polypeptides of the invention can be combined with domains derived from other transcription factors or synthetic domains to modify the biological activity of a transcription factor. For instance, a DNA-binding domain derived from a transcription factor of the invention can be combined with the activation domain of another transcription
15 factor or with a synthetic activation domain. A transcription activation domain assists in initiating transcription from a DNA-binding site. Examples include the transcription activation region of VP16 or GAL4 (Moore et al. (1998) *Proc. Natl. Acad. Sci.* 95: 376-381; Aoyama et al. (1995) *Plant Cell* 7: 1773-1785), peptides derived from bacterial sequences (Ma and Ptashne (1987) *Cell* 51: 113-119) and synthetic peptides (Giniger and Ptashne (1987) *Nature* 330: 670-672).

Expression and Modification of Polypeptides

Typically, polynucleotide sequences of the invention are incorporated into recombinant DNA (or RNA) molecules that direct expression of polypeptides of the invention in appropriate host cells, transgenic plants, in vitro translation systems, or the like. Due to the inherent degeneracy of the genetic
25 code, nucleic acid sequences which encode substantially the same or a functionally equivalent amino acid sequence can be substituted for any listed sequence to provide for cloning and expressing the relevant homolog.

The transgenic plants of the present invention comprising recombinant polynucleotide sequences are generally derived from parental plants, which may themselves be non-transformed (or non-transgenic)
30 plants. These transgenic plants may either have a transcription factor gene "knocked out" (for example, with a genomic insertion by homologous recombination, an antisense or ribozyme construct) or expressed to a normal or wild-type extent. However, overexpressing transgenic "progeny" plants will exhibit greater mRNA levels, wherein the mRNA encodes a transcription factor, that is, a DNA-binding protein that is

capable of binding to a DNA regulatory sequence and inducing transcription, and preferably, expression of a plant trait gene. Preferably, the mRNA expression level will be at least three-fold greater than that of the parental plant, or more preferably at least ten-fold greater mRNA levels compared to said parental plant, and most preferably at least fifty-fold greater compared to said parental plant.

5

Vectors, Promoters, and Expression Systems

The present invention includes recombinant constructs comprising one or more of the nucleic acid sequences herein. The constructs typically comprise a vector, such as a plasmid, a cosmid, a phage, a virus (e.g., a plant virus), a bacterial artificial chromosome (BAC), a yeast artificial chromosome (YAC), or the like, into which a nucleic acid sequence of the invention has been inserted, in a forward or reverse orientation. In a preferred aspect of this embodiment, the construct further comprises regulatory sequences, including, for example, a promoter, operably linked to the sequence. Large numbers of suitable vectors and promoters are known to those of skill in the art, and are commercially available.

General texts that describe molecular biological techniques useful herein, including the use and production of vectors, promoters and many other relevant topics, include Berger, Sambrook, *supra*, and Ausubel, *supra*. Any of the identified sequences can be incorporated into a cassette or vector, e.g., for expression in plants. A number of expression vectors suitable for stable transformation of plant cells or for the establishment of transgenic plants have been described including those described in Weissbach and Weissbach (1989) Methods for Plant Molecular Biology, Academic Press, and Gelvin et al. (1990) Plant Molecular Biology Manual, Kluwer Academic Publishers. Specific examples include those derived from a Ti plasmid of *Agrobacterium tumefaciens*, as well as those disclosed by Herrera-Estrella et al. (1983) *Nature* 303: 209, Bevan (1984) *Nucleic Acids Res.* 12: 8711-8721, Klee (1985) *Bio/Technology* 3: 637-642, for dicotyledonous plants.

Alternatively, non-Ti vectors can be used to transfer the DNA into monocotyledonous plants and cells by using free DNA delivery techniques. Such methods can involve, for example, the use of liposomes, electroporation, microprojectile bombardment, silicon carbide whiskers, and viruses. By using these methods transgenic plants such as wheat, rice (Christou (1991) *Bio/Technology* 9: 957-962) and corn (Gordon-Kamm (1990) *Plant Cell* 2: 603-618) can be produced. An immature embryo can also be a good target tissue for monocots for direct DNA delivery techniques by using the particle gun (Weeks et al. (1993) *Plant Physiol.* 102: 1077-1084; Vasil (1993) *Bio/Technology* 10: 667-674; Wan and Lemeaux (1994) *Plant Physiol.* 104: 37-48, and for *Agrobacterium*-mediated DNA transfer (Ishida et al. (1996) *Nature Biotechnol.* 14: 745-750).

Typically, plant transformation vectors include one or more cloned plant coding sequence

(genomic or cDNA) under the transcriptional control of 5' and 3' regulatory sequences and a dominant selectable marker. Such plant transformation vectors typically also contain a promoter (e.g., a regulatory region controlling inducible or constitutive, environmentally-or developmentally-regulated, or cell- or tissue-specific expression), a transcription initiation start site, an RNA processing signal (such as intron splice sites), a transcription termination site, and/or a polyadenylation signal.

A potential utility for the transcription factor polynucleotides disclosed herein is the isolation of promoter elements from these genes that can be used to program expression in plants of any genes. Each transcription factor gene disclosed herein is expressed in a unique fashion, as determined by promoter elements located upstream of the start of translation, and additionally within an intron of the transcription factor gene or downstream of the termination codon of the gene. As is well known in the art, for a significant portion of genes, the promoter sequences are located entirely in the region directly upstream of the start of translation. In such cases, typically the promoter sequences are located within 2.0 kb of the start of translation, or within 1.5 kb of the start of translation, frequently within 1.0 kb of the start of translation, and sometimes within 0.5 kb of the start of translation.

The promoter sequences can be isolated according to methods known to one skilled in the art.

Examples of constitutive plant promoters which can be useful for expressing the TF sequence include: the cauliflower mosaic virus (CaMV) 35S promoter, which confers constitutive, high-level expression in most plant tissues (see, for example, Odell et al. (1985) *Nature* 313: 810-812); the nopaline synthase promoter (An et al. (1988) *Plant Physiol.* 88: 547-552); and the octopine synthase promoter (Fromm et al. (1989) *Plant Cell* 1: 977-984).

The transcription factors of the invention may be operably linked with a specific promoter that causes the transcription factor to be expressed in response to environmental, tissue-specific or temporal signals. A variety of plant gene promoters are known to regulate gene expression in response to environmental, hormonal, chemical, developmental signals, and in a tissue-active manner; many of these may be used for expression of a TF sequence in plants. Choice of a promoter is based largely on the phenotype of interest and is determined by such factors as tissue (e.g., seed, fruit, root, pollen, vascular tissue, flower, carpel, etc.), inducibility (e.g., in response to wounding, heat, cold, drought, light, pathogens, etc.), timing, developmental stage, and the like. Numerous known promoters have been characterized and can favorably be employed to promote expression of a polynucleotide of the invention in a transgenic plant or cell of interest. For example, tissue specific promoters include: seed-specific promoters (such as the napin, phaseolin or DC3 promoter described in US Pat. No. 5,773,697), fruit-specific promoters that are active during fruit ripening (such as the *dru 1* promoter (US Pat. No. 5,783,393), or the *2A11* promoter (US Pat. No. 4,943,674) and the tomato polygalacturonase promoter

(Bird et al. (1988) *Plant Mol. Biol.* 11: 651-662), root-specific promoters, such as ARSK1, and those disclosed in US Patent Nos. 5,618,988, 5,837,848 and 5,905,186, epidermis-specific promoters, including CUT1 (Kunst et al. (1999) *Biochem. Soc. Trans.* 28: 651-654), pollen-active promoters such as *PTA29*, *PTA26* and *PTA13* (US Pat. No. 5,792,929), promoters active in vascular tissue (Ringli and Keller (1998) *Plant Mol. Biol.* 37: 977-988), flower-specific (Kaiser et al. (1995) *Plant Mol. Biol.* 28: 231-243), pollen (Baerson et al. (1994) *Plant Mol. Biol.* 26: 1947-1959), carpels (Ohl et al. (1990) *Plant Cell* 2: 837-848), pollen and ovules (Baerson et al. (1993) *Plant Mol. Biol.* 22: 255-267), auxin-inducible promoters (such as that described in van der Kop et al. (1999) *Plant Mol. Biol.* 39: 979-990 or Baumann et al. (1999) *Plant Cell* 11: 323-334), cytokinin-inducible promoter (Guevara-Garcia (1998) *Plant Mol. Biol.* 38: 743-753), promoters responsive to gibberellin (Shi et al. (1998) *Plant Mol. Biol.* 38: 1053-1060, Willmott et al. (1998) *Plant Mol. Biol.* 38: 817-825) and the like. Additional promoters are those that elicit expression in response to heat (Ainley et al. (1993) *Plant Mol. Biol.* 22: 13-23), light (e.g., the pea *rbcS*-3A promoter, Kuhlmeier et al. (1989) *Plant Cell* 1: 471-478, and the maize *rbcS* promoter, Schaffner and Sheen (1991) *Plant Cell* 3: 997-1012); wounding (e.g., *wun1*, Siebertz et al. (1989) *Plant Cell* 1: 961-968); pathogens (such as the PR-1 promoter described in Buchel et al. (1999) *Plant Mol. Biol.* 40: 387-396, and the PDF1.2 promoter described in Manners et al. (1998) *Plant Mol. Biol.* 38: 1071-1080), and chemicals such as methyl jasmonate or salicylic acid (Gatz (1997) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* 48: 89-108). In addition, the timing of the expression can be controlled by using promoters such as those acting at senescence (Gan and Amasino (1995) *Science* 270: 1986-1988); or late seed development (Odell et al. (1994) *Plant Physiol.* 106: 447-458).

Plant expression vectors can also include RNA processing signals that can be positioned within, upstream or downstream of the coding sequence. In addition, the expression vectors can include additional regulatory sequences from the 3'-untranslated region of plant genes, e.g., a 3' terminator region to increase mRNA stability of the mRNA, such as the PI-II terminator region of potato or the octopine or nopaline synthase 3' terminator regions.

Additional Expression Elements

Specific initiation signals can aid in efficient translation of coding sequences. These signals can include, e.g., the ATG initiation codon and adjacent sequences. In cases where a coding sequence, its initiation codon and upstream sequences are inserted into the appropriate expression vector, no additional translational control signals may be needed. However, in cases where only coding sequence (e.g., a mature protein coding sequence), or a portion thereof, is inserted, exogenous transcriptional control signals including the ATG initiation codon can be separately provided. The initiation codon is provided in the

correct reading frame to facilitate transcription. Exogenous transcriptional elements and initiation codons can be of various origins, both natural and synthetic. The efficiency of expression can be enhanced by the inclusion of enhancers appropriate to the cell system in use.

5 Expression Hosts

The present invention also relates to host cells which are transduced with vectors of the invention, and the production of polypeptides of the invention (including fragments thereof) by recombinant techniques. Host cells are genetically engineered (i.e., nucleic acids are introduced, e.g., transduced, transformed or transfected) with the vectors of this invention, which may be, for example, a cloning vector
10 or an expression vector comprising the relevant nucleic acids herein. The vector is optionally a plasmid, a viral particle, a phage, a naked nucleic acid, etc. The engineered host cells can be cultured in conventional nutrient media modified as appropriate for activating promoters, selecting transformants, or amplifying the relevant gene. The culture conditions, such as temperature, pH and the like, are those previously used with the host cell selected for expression, and will be apparent to those skilled in the art and in the references
15 cited herein, including, Sambrook, *supra* and Ausubel, *supra*.

The host cell can be a eukaryotic cell, such as a yeast cell, or a plant cell, or the host cell can be a prokaryotic cell, such as a bacterial cell. Plant protoplasts are also suitable for some applications. For example, the DNA fragments are introduced into plant tissues, cultured plant cells or plant protoplasts by standard methods including electroporation (Fromm et al. (1985) *Proc. Natl. Acad. Sci.* 82: 5824-5828,
20 infection by viral vectors such as cauliflower mosaic virus (CaMV) (Hohn et al. (1982) Molecular Biology of Plant Tumors, Academic Press, New York, NY, pp. 549-560; US 4,407,956), high velocity ballistic penetration by small particles with the nucleic acid either within the matrix of small beads or particles, or on the surface (Klein et al. (1987) *Nature* 327: 70-73), use of pollen as vector (WO 85/01856), or use of *Agrobacterium tumefaciens* or *A. rhizogenes* carrying a T-DNA plasmid in which DNA fragments are
25 cloned. The T-DNA plasmid is transmitted to plant cells upon infection by *Agrobacterium tumefaciens*, and a portion is stably integrated into the plant genome (Horsch et al. (1984) *Science* 233: 496-498; Fraley et al. (1983) *Proc. Natl. Acad. Sci.* 80: 4803-4807).

The cell can include a nucleic acid of the invention that encodes a polypeptide, wherein the cell expresses a polypeptide of the invention. The cell can also include vector sequences, or the like.
30 Furthermore, cells and transgenic plants that include any polypeptide or nucleic acid above or throughout this specification, e.g., produced by transduction of a vector of the invention, are an additional feature of the invention.

For long-term, high-yield production of recombinant proteins, stable expression can be used. Host

cells transformed with a nucleotide sequence encoding a polypeptide of the invention are optionally cultured under conditions suitable for the expression and recovery of the encoded protein from cell culture. The protein or fragment thereof produced by a recombinant cell may be secreted, membrane-bound, or contained intracellularly, depending on the sequence and/or the vector used. As will be understood by those of skill in the art, expression vectors containing polynucleotides encoding mature proteins of the invention can be designed with signal sequences which direct secretion of the mature polypeptides through a prokaryotic or eukaryotic cell membrane.

Modified Amino Acid Residues

Polypeptides of the invention may contain one or more modified amino acid residues. The presence of modified amino acids may be advantageous in, for example, increasing polypeptide half-life, reducing polypeptide antigenicity or toxicity, increasing polypeptide storage stability, or the like. Amino acid residue(s) are modified, for example, co-translationally or post-translationally during recombinant production or modified by synthetic or chemical means.

Non-limiting examples of a modified amino acid residue include incorporation or other use of acetylated amino acids, glycosylated amino acids, sulfated amino acids, prenylated (e.g., farnesylated, geranylgeranylated) amino acids, PEG modified (for example, "PEGylated") amino acids, biotinylated amino acids, carboxylated amino acids, phosphorylated amino acids, etc. References adequate to guide one of skill in the modification of amino acid residues are replete throughout the literature.

The modified amino acid residues may prevent or increase affinity of the polypeptide for another molecule, including, but not limited to, polynucleotide, proteins, carbohydrates, lipids and lipid derivatives, and other organic or synthetic compounds.

Identification of Additional Protein Factors

A transcription factor provided by the present invention can also be used to identify additional endogenous or exogenous molecules that can affect a phenotype or trait of interest. Such molecules include endogenous molecules that are acted upon either at a transcriptional level by a transcription factor of the invention to modify a phenotype as desired. For example, the transcription factors can be employed to identify one or more downstream genes that are subject to a regulatory effect of the transcription factor.

In one approach, a transcription factor or transcription factor homolog of the invention is expressed in a host cell, e.g., a transgenic plant cell, tissue or explant, and expression products, either RNA or protein, of likely or random targets are monitored, e.g., by hybridization to a microarray of nucleic acid probes corresponding to genes expressed in a tissue or cell type of interest, by two-dimensional gel

electrophoresis of protein products, or by any other method known in the art for assessing expression of gene products at the level of RNA or protein. Alternatively, a transcription factor of the invention can be used to identify promoter sequences (such as binding sites on DNA sequences) involved in the regulation of a downstream target. After identifying a promoter sequence, interactions between the transcription factor and the promoter sequence can be modified by changing specific nucleotides in the promoter sequence or specific amino acids in the transcription factor that interact with the promoter sequence to alter a plant trait. Typically, transcription factor DNA-binding sites are identified by gel shift assays. After identifying the promoter regions, the promoter region sequences can be employed in double-stranded DNA arrays to identify molecules that affect the interactions of the transcription factors with their promoters (Bulyk et al. (1999) *Nature Biotechnol.* 17: 573-577).

The identified transcription factors are also useful to identify proteins that modify the activity of the transcription factor. Such modification can occur by covalent modification, such as by phosphorylation, or by protein-protein (homo or-heteropolymer) interactions. Any method suitable for detecting protein-protein interactions can be employed. Among the methods that can be employed are co-immunoprecipitation, cross-linking and co-purification through gradients or chromatographic columns, and the two-hybrid yeast system.

The two-hybrid system detects protein interactions *in vivo* and has been previously described (Chien et al. (1991) *Proc. Natl. Acad. Sci.* 88: 9578-9582), and is commercially available from Clontech (Palo Alto, Calif.). In such a system, plasmids are constructed that encode two hybrid proteins: one consists of the DNA-binding domain of a transcription activator protein fused to the TF polypeptide and the other consists of the transcription activator protein's activation domain fused to an unknown protein that is encoded by a cDNA that has been recombined into the plasmid as part of a cDNA library. The DNA-binding domain fusion plasmid and the cDNA library are transformed into a strain of the yeast *Saccharomyces cerevisiae* that contains a reporter gene (e.g., lacZ) whose regulatory region contains the transcription activator's binding site. Either hybrid protein alone cannot activate transcription of the reporter gene. Interaction of the two hybrid proteins reconstitutes the functional activator protein and results in expression of the reporter gene, which is detected by an assay for the reporter gene product. Then, the library plasmids responsible for reporter gene expression are isolated and sequenced to identify the proteins encoded by the library plasmids. After identifying proteins that interact with the transcription factors, assays for compounds that interfere with the TF protein-protein interactions can be preformed.

Subsequences

Also contemplated are uses of polynucleotides, also referred to herein as oligonucleotides, typically having at least 12 bases, preferably at least 15, more preferably at least 20, 30, or 50 bases, which hybridize under at least highly stringent (or ultra-high stringent or ultra-ultra-high stringent conditions) conditions to a polynucleotide sequence described above. The polynucleotides may be used as probes, primers, sense and antisense agents, and the like, according to methods as noted *supra*.

Subsequences of the polynucleotides of the invention, including polynucleotide fragments and oligonucleotides are useful as nucleic acid probes and primers. An oligonucleotide suitable for use as a probe or primer is at least about 15 nucleotides in length, more often at least about 18 nucleotides, often at least about 21 nucleotides, frequently at least about 30 nucleotides, or about 40 nucleotides, or more in length. A nucleic acid probe is useful in hybridization protocols, for example, to identify additional polypeptide homologs of the invention, including protocols for microarray experiments. Primers can be annealed to a complementary target DNA strand by nucleic acid hybridization to form a hybrid between the primer and the target DNA strand, and then extended along the target DNA strand by a DNA polymerase enzyme. Primer pairs can be used for amplification of a nucleic acid sequence, e.g., by the polymerase chain reaction (PCR) or other nucleic-acid amplification methods. See Sambrook, *supra*, and Ausubel, *supra*.

In addition, the invention includes an isolated or recombinant polypeptide including a subsequence of at least about 15 contiguous amino acids encoded by the recombinant or isolated polynucleotides of the invention. For example, such polypeptides, or domains or fragments thereof, can be used as immunogens, e.g., to produce antibodies specific for the polypeptide sequence, or as probes for detecting a sequence of interest. A subsequence can range in size from about 15 amino acids in length up to and including the full length of the polypeptide.

To be encompassed by the present invention, an expressed polypeptide which comprises such a polypeptide subsequence performs at least one biological function of the intact polypeptide in substantially the same manner, or to a similar extent, as does the intact polypeptide. For example, a polypeptide fragment can comprise a recognizable structural motif or functional domain such as a DNA binding domain that activates transcription, for example, by binding to a specific DNA promoter region an activation domain, or a domain for protein-protein interactions.

Production of Transgenic Plants

Modification of Traits

The polynucleotides of the invention are favorably employed to produce transgenic plants with

various traits, or characteristics, that have been modified in a desirable manner, e.g., to improve the seed characteristics of a plant. For example, alteration of expression levels or patterns (e.g., spatial or temporal expression patterns) of one or more of the transcription factors (or transcription factor homologs) of the invention, as compared with the levels of the same protein found in a wild-type plant, can be used to modify a plant's traits. An illustrative example of trait modification, improved characteristics, by altering expression levels of a particular transcription factor is described further in the Examples and the Sequence Listing.

Arabidopsis as a model system

Arabidopsis thaliana is the object of rapidly growing attention as a model for genetics and metabolism in plants. *Arabidopsis* has a small genome, and well-documented studies are available. It is easy to grow in large numbers and mutants defining important genetically controlled mechanisms are either available, or can readily be obtained. Various methods to introduce and express isolated homologous genes are available (see Koncz et al., editors, Methods in Arabidopsis Research (1992) World Scientific, New Jersey NJ, in "Preface"). Because of its small size, short life cycle, obligate autogamy and high fertility, *Arabidopsis* is also a choice organism for the isolation of mutants and studies in morphogenetic and development pathways, and control of these pathways by transcription factors (Koncz *supra*, p. 72). A number of studies introducing transcription factors into *A. thaliana* have demonstrated the utility of this plant for understanding the mechanisms of gene regulation and trait alteration in plants. (See, for example, Koncz *supra*, and U.S. Patent Number 6,417,428).

Arabidopsis genes in transgenic plants

Expression of genes which encode transcription factors modify expression of endogenous genes, polynucleotides, and proteins are well known in the art. In addition, transgenic plants comprising isolated polynucleotides encoding transcription factors may also modify expression of endogenous genes, polynucleotides, and proteins. Examples include Peng et al. (1997) et al. *Genes and Development* 11: 3194-3205, and Peng et al. (1999) *Nature* 400: 256-261. In addition, many others have demonstrated that an *Arabidopsis* transcription factor expressed in an exogenous plant species elicits the same or very similar phenotypic response. See, for example, Fu et al. (2001) *Plant Cell* 13: 1791-1802; Nandi et al. (2000) *Curr. Biol.* 10: 215-218; Coupland (1995) *Nature* 377: 482-483; and Weigel and Nilsson (1995) *Nature* 377: 482-500.

Homologous genes introduced into transgenic plants

Homologous genes that may be derived from any plant, or from any source whether natural, synthetic, semi-synthetic or recombinant, and that share significant sequence identity or similarity to those provided by the present invention, may be introduced into plants, for example, crop plants, to confer desirable or improved traits. Consequently, transgenic plants may be produced that comprise a recombinant expression vector or cassette with a promoter operably linked to one or more sequences homologous to presently disclosed sequences. The promoter may be, for example, a plant or viral promoter.

The invention thus provides for methods for preparing transgenic plants, and for modifying plant traits. These methods include introducing into a plant a recombinant expression vector or cassette comprising a functional promoter operably linked to one or more sequences homologous to presently disclosed sequences. Plants and kits for producing these plants that result from the application of these methods are also encompassed by the present invention.

Transcription factors of interest for the modification of plant traits

Currently, the existence of a series of maturity groups for different latitudes represents a major barrier to the introduction of new valuable traits. Any trait (e.g. abiotic stress tolerance or increased biomass) has to be bred into each of the different maturity groups separately, a laborious and costly exercise. The availability of single strain, which could be grown at any latitude, would therefore greatly increase the potential for introducing new traits to crop species such as soybean and cotton.

For the specific effects, traits and utilities conferred to plants, one or more transcription factor genes of the present invention may be used to increase or decrease, or improve or prove deleterious to a given trait. For example, knocking out a transcription factor gene that naturally occurs in a plant, or suppressing the gene (with, for example, antisense suppression), may cause decreased tolerance to an osmotic stress relative to non-transformed or wild-type plants. By overexpressing this gene, the plant may experience increased tolerance to the same stress. More than one transcription factor gene may be introduced into a plant, either by transforming the plant with one or more vectors comprising two or more transcription factors, or by selective breeding of plants to yield hybrid crosses that comprise more than one introduced transcription factor.

Genes, traits and utilities that affect plant characteristics

Plant transcription factors can modulate gene expression, and, in turn, be modulated by the environmental experience of a plant. Significant alterations in a plant's environment invariably result in a

change in the plant's transcription factor gene expression pattern. Altered transcription factor expression patterns generally result in phenotypic changes in the plant. Transcription factor gene product(s) in transgenic plants then differ(s) in amounts or proportions from that found in wild-type or non-transformed plants, and those transcription factors likely represent polypeptides that are used to alter the response to the environmental change. By way of example, it is well accepted in the art that analytical methods based on altered expression patterns may be used to screen for phenotypic changes in a plant far more effectively than can be achieved using traditional methods.

Increased biomass.

Plants overexpressing nine distinct related AT-hook transcription factors of the invention, including sequences from diverse species of monocots and dicots, such as *Arabidopsis thaliana* polypeptides G1073, G1067, G2153 and G2156, *Oryza sativa* polypeptides G3399 and G3407, and *Glycine max* polypeptides G3456, G3459 and G3460, become larger than controls, and generally produce broader leaves than wild-type plants. For some ornamental plants, the ability to provide larger varieties with these genes or their equivalents may be highly desirable. More significantly, crop species overexpressing these genes from diverse species would also produce higher yields on larger cultivars, particularly those in which the vegetative portion of the plant is edible. This has already been observed in *Arabidopsis* and tomato plants. Tomato plants overexpressing the *A. thaliana* G2153 polypeptide have been found to be larger and produce more fruit than wild-type control tomato plants. Numerous *Arabidopsis* lines that overexpress G3399 and G3407, which are rice genes, and G3456, G3459 and G3460, which are soy genes, develop significantly larger rosettes and leaves than wild-type *Arabidopsis* controls.

Overexpression of these genes can confer increased stress tolerance as well as increased biomass, and the increased biomass appears to be related to the particular mechanism of stress tolerance exhibited by these genes. The decision for a lateral organ to continue growth and expansion versus entering late development phases (growth cessation and senescence) is controlled genetically and hormonally, including regulation at an organ size checkpoint (e.g., Mizukami (2001) *Curr Opinion Plant Biol* 4: 533-39; Mizukami and Fisher (2000) *Proc. Natl. Acad. Sci.* 97: 942-47; Hu et al. 2003 *Plant Cell* 15: 1591). Organ size is controlled by the meristematic competence of organ cells, with increased meristematic competence leading to increased organ size (both leaves and stems). Plant hormones can impact plant organ size, with, for example, ethylene pathway overexpression leading to reduced organ size. There also suggestions that auxin plays a determinative role in organ size. Stress responses can impact hormone levels in plant tissues, including ABA and ethylene levels, thereby modifying meristematic competence and final organ size. Thus, overexpression of *HRC* genes alters environmental (e.g., stress) inputs to the

organ size checkpoint, thus enhancing organ size under typical growth conditions.

Due to frequent exposure to stresses under typical plant growth conditions, the maximum genetically programmed organ size is infrequently achieved. It is well appreciated that increased leaf organ size can result in increased seed yield, through enhanced energy capture and source activity. Thus, a major strategy for yield optimization is altered characteristics of the sensor that integrates external environmental stress inputs to meristematic competence and organ size control. The HRC genes that are the subject of the instant invention represent one component of this control mechanism. Increased expression of HRC genes leads to diminished sensitivity of the environmental sensor for organ size control to those stress inputs. This increase in stress threshold for diminished meristematic competence results in increased vegetative and seed yield under typical plant growth conditions. AT-hook proteins are known to modulate gene expression through interactions with other proteins. Thus, the environmental integration mechanism for organ size control instantiated by HRC proteins will have additional components whose function will be recognized by the ability of the encoded proteins to participate in regulating gene sets that are regulated by HRC proteins. Identification of additional components of the integration can be achieved by identifying other transcription factors that bind to upstream regulatory regions, detecting proteins that directly interact with HRC proteins.

Sugar sensing.

In addition to their important role as an energy source and structural component of the plant cell, sugars are central regulatory molecules that control several aspects of plant physiology, metabolism and development (Hsieh et al. (1998) *Proc. Natl. Acad. Sci.* 95: 13965-13970). It is thought that this control is achieved by regulating gene expression and, in higher plants, sugars have been shown to repress or activate plant genes involved in many essential processes such as photosynthesis, glyoxylate metabolism, respiration, starch and sucrose synthesis and degradation, pathogen response, wounding response, cell cycle regulation, pigmentation, flowering and senescence. The mechanisms by which sugars control gene expression are not understood.

Several sugar sensing mutants have turned out to be allelic to abscisic acid (ABA) and ethylene mutants. ABA is found in all photosynthetic organisms and acts as a key regulator of transpiration, stress responses, embryogenesis, and seed germination. Most ABA effects are related to the compound acting as a signal of decreased water availability, whereby it triggers a reduction in water loss, slows growth, and mediates adaptive responses. However, ABA also influences plant growth and development via interactions with other phytohormones. Physiological and molecular studies indicate that maize and *Arabidopsis* have almost identical pathways with regard to ABA biosynthesis and signal transduction. For

further review, see Finkelstein and Rock ((2002) Absciscic acid biosynthesis and response (In The Arabidopsis Book, Editors: Somerville and Meyerowitz (American Society of Plant Biologists, Rockville, MD).

This potentially implicates G1073, G2153, G2156 and related transcription factors in hormone signaling based on the sucrose sugar sensing phenotype of 35S::G1073, 35S::G2153 and 35S::G2156 transgenic lines. On the other hand, the sucrose treatment used in these experiments (9.5% w/v) could also be an osmotic stress. Therefore, one could interpret these data as an indication that the 35S::G1073, 35S::G2153 and 35S::G2156 transgenic lines are more tolerant to osmotic stress. However, it is well known that plant responses to ABA, osmotic and other stress may be linked, and these different treatments may even act in a synergistic manner to increase the degree of a response. For example, Xiong, Ishitani, and Zhu ((1999) *Plant Physiol.* 119: 205-212) have shown that genetic and molecular studies may be used to show extensive interaction between osmotic stress, temperature stress, and ABA responses in plants. These investigators analyzed the expression of *RD29A-LUC* in response to various treatment regimes in *Arabidopsis*. The RD29A promoter contains both the ABA-responsive and the dehydration-responsive element - also termed the C-repeat - and can be activated by osmotic stress, low temperature, or ABA treatment; transcription of the RD29A gene in response to osmotic and cold stresses is mediated by both ABA-dependent and ABA-independent pathways (Xiong, Ishitani, and Zhu (1999) *supra*). LUC refers to the firefly luciferase coding sequence, which, in this case, was driven by the stress responsive RD29A promoter. The results revealed both positive and negative interactions, depending on the nature and duration of the treatments. Low temperature stress was found to impair osmotic signaling but moderate heat stress strongly enhanced osmotic stress induction, thus acting synergistically with osmotic signaling pathways. In this study, the authors reported that osmotic stress and ABA can act synergistically by showing that the treatments simultaneously induced transgene and endogenous gene expression. Similar results were reported by Bostock and Quatrano ((1992) *Plant Physiol.* 98: 1356-1363), who found that osmotic stress and ABA act synergistically and induce maize *Em* gene expression. Ishitani et al (1997) *Plant Cell* 9: 1935-1949) isolated a group of *Arabidopsis* single-gene mutations that confer enhanced responses to both osmotic stress and ABA. The nature of the recovery of these mutants from osmotic stress and ABA treatment suggested that although separate signaling pathways exist for osmotic stress and ABA, the pathways share a number of components; these common components may mediate synergistic interactions between osmotic stress and ABA. Thus, contrary to the previously-held belief that ABA-dependent and ABA-independent stress signaling pathways act in a parallel manner, our data reveal that these pathways cross-talk and converge to activate stress gene expression.

Because sugars are important signaling molecules, the ability to control either the concentration of

a signaling sugar or how the plant perceives or responds to a signaling sugar could be used to control plant development, physiology or metabolism. For example, the flux of sucrose (a disaccharide sugar used for systemically transporting carbon and energy in most plants) has been shown to affect gene expression and alter storage compound accumulation in seeds. Manipulation of the sucrose signaling pathway in seeds may therefore cause seeds to have more protein, oil or carbohydrate, depending on the type of manipulation. Similarly, in tubers, sucrose is converted to starch which is used as an energy store. It is thought that sugar signaling pathways may partially determine the levels of starch synthesized in the tubers. The manipulation of sugar signaling in tubers could lead to tubers with a higher starch content.

Thus, the presently disclosed transcription factor genes that manipulate the sugar signal transduction pathway, including, for example, G1073 and G2156, along with their equivalents, may lead to altered gene expression to produce plants with desirable traits. In particular, manipulation of sugar signal transduction pathways could be used to alter source-sink relationships in seeds, tubers, roots and other storage organs leading to increase in yield.

Salt and drought tolerance

Plants are subject to a range of environmental challenges. Several of these, including salt stress, general osmotic stress, drought stress and freezing stress, have the ability to impact whole plant and cellular water availability. Not surprisingly, then, plant responses to this collection of stresses are related. In a recent review, Zhu notes that “most studies on water stress signaling have focused on salt stress primarily because plant responses to salt and drought are closely related and the mechanisms overlap” (Zhu (2002) *Ann. Rev. Plant Biol.* 53: 247-273). Many examples of similar responses (i.e., genetic pathways to this set of stresses have been documented. For example, the CBF transcription factors have been shown to condition resistance to salt, freezing and drought (Kasuga et al. (1999) *Nature Biotech.* 17: 287-291). The *Arabidopsis rd29B* gene is induced in response to both salt and dehydration stress, a process that is mediated largely through an ABA signal transduction process (Uno et al. (2000) *Proc. Natl. Acad. Sci. USA* 97: 11632-11637), resulting in altered activity of transcription factors that bind to an upstream element within the *rd29B* promoter. In *Mesembryanthemum crystallinum* (ice plant), Patharker and Cushman have shown that a calcium-dependent protein kinase (McCDPK1) is induced by exposure to both drought and salt stresses (Patharker and Cushman (2000) *Plant J.* 24: 679-691). The stress-induced kinase was also shown to phosphorylate a transcription factor, presumably altering its activity, although transcript levels of the target transcription factor are not altered in response to salt or drought stress. Similarly, Saijo et al. demonstrated that a rice salt/drought-induced calmodulin-dependent protein kinase (OsCDPK7) conferred increased salt and drought tolerance to rice when overexpressed (Saijo et al. (2000) *Plant J.* 23: 319-327).

Exposure to dehydration invokes similar survival strategies in plants as does freezing stress (see, for example, Yelenosky (1989) *Plant Physiol* 89: 444-451) and drought stress induces freezing tolerance (see, for example, Siminovitch et al. (1982) *Plant Physiol* 69: 250-255; and Guy et al. (1992) *Planta* 188: 265-270). In addition to the induction of cold-acclimation proteins, strategies that allow plants to survive in low water conditions may include, for example, reduced surface area, or surface oil or wax production. Plants overexpressing G1073, G1067 and G2156 have been shown to be more tolerant to drought stress than wild-type control plants.

Consequently, one skilled in the art would expect that some pathways involved in resistance to one of these stresses, and hence regulated by an individual transcription factor, will also be involved in resistance to another of these stresses, regulated by the same or homologous transcription factors. Of course, the overall resistance pathways are related, not identical, and therefore not all transcription factors controlling resistance to one stress will control resistance to the other stresses. Nonetheless, if a transcription factor conditions resistance to one of these stresses, it would be apparent to one skilled in the art to test for resistance to these related stresses.

Thus, modifying the expression of a number of presently disclosed transcription factor genes, including G1073, G1067 and G2156 and their equivalents, may be used to increase a plant's tolerance to low water conditions and provide the benefits of improved survival, increased yield and an extended geographic and temporal planting range.

Osmotic stress. A number of these genes (G1073, G1067, G2153 and G2156) have been shown to have an altered osmotic stress tolerance phenotype, by virtue of their improved germination on high sugar-containing media. Most of these genes have also been shown to confer increased salt stress and drought tolerance to overexpressing plants (all have been shown to increase osmotic stress tolerance in *Arabidopsis*, and G2153 has been shown to do the same in tomatoes). Thus, modification of the expression of these and other structurally related disclosed transcription factor genes may be used to increase germination rate or growth under adverse osmotic conditions, which could impact survival and yield of seeds and plants. Osmotic stresses may be regulated by specific molecular control mechanisms that include genes controlling water and ion movements, functional and structural stress-induced proteins, signal perception and transduction, and free radical scavenging, and many others (Wang et al. (2001) *Acta Hort.* (ISHS) 560: 285-292). Instigators of osmotic stress include freezing, drought and high salinity, each of which are discussed in more detail below.

In many ways, freezing, high salt and drought have similar effects on plants, not the least of which is the induction of common polypeptides that respond to these different stresses. For example, freezing is

similar to water deficit in that freezing reduces the amount of water available to a plant. Exposure to freezing temperatures may lead to cellular dehydration as water leaves cells and forms ice crystals in intercellular spaces (Buchanan, *supra*). As with high salt concentration and freezing, the problems for plants caused by low water availability include mechanical stresses caused by the withdrawal of cellular water. Thus, the incorporation of transcription factors that modify a plant's response to osmotic stress into, for example, a crop or ornamental plant, may be useful in reducing damage or loss. Specific effects caused by freezing, high salt and drought are addressed below.

Salt. The genes of the Sequence Listing, including, for example, G1073, G1067 and G2156, that provide tolerance to salt may be used to engineer salt tolerant crops and trees that can flourish in soils with high saline content or under drought conditions. In particular, increased salt tolerance during the germination stage of a plant enhances survival and yield. Presently disclosed transcription factor genes that provide increased salt tolerance during germination, the seedling stage, and throughout a plant's life cycle, would find particular value for imparting survival and yield in areas where a particular crop would not normally prosper.

Summary of altered plant characteristics. A clade of structurally and functionally related sequences that derive from a wide range of plants, including polynucleotide *Arabidopsis* SEQ ID NOs: 1, 3, 5, 7, fragments thereof, rice SEQ ID NOs: 9, 11, and soy SEQ ID NOs: 13, 15, and 17, fragments thereof, paralogs, orthologs, equivalent, and fragments thereof, is provided. These sequences have been shown in laboratory and field experiments to confer altered size and abiotic stress tolerance phenotypes in plants. The invention also provides polypeptides comprising: *Arabidopsis* SEQ ID NOs: 2, 4, 6, 8, rice SEQ ID NOs: 10, 12, and soy SEQ ID NOs: 14, 16, 18, and fragments thereof, conserved domains thereof, paralogs, orthologs, equivalent, and fragments thereof. Plants that overexpress these sequences have been observed to become larger, and a significant number have been shown to be more tolerant to a wide variety of abiotic stresses, including, for example, osmotic stresses such as drought and high salt levels. Many of the orthologs of these sequences are listed in the Sequence Listing, and due to the high degree of structural similarity to the sequences of the invention, it is expected that these sequences may also function to increase plant biomass and/or abiotic stress tolerance. The invention also encompasses the complements of the polynucleotides. The polynucleotides are useful for screening libraries of molecules or compounds for specific binding and for creating transgenic plants having increased biomass and/or abiotic stress tolerance.

Antisense and Co-suppression

In addition to expression of the nucleic acids of the invention as gene replacement or plant phenotype modification nucleic acids, the nucleic acids are also useful for sense and anti-sense suppression of expression, e.g., to down-regulate expression of a nucleic acid of the invention, e.g., as a further mechanism for modulating plant phenotype. That is, the nucleic acids of the invention, or subsequences or anti-sense sequences thereof, can be used to block expression of naturally occurring homologous nucleic acids. A variety of sense and anti-sense technologies are known in the art, e.g., as set forth in Lichtenstein and Nellen (1997) Antisense Technology: A Practical Approach IRL Press at Oxford University Press, Oxford, U.K. Antisense regulation is also described in Crowley et al. (1985) *Cell* 43: 633-641; Rosenberg et al. (1985) *Nature* 313: 703-706; Preiss et al. (1985) *Nature* 313: 27-32; Melton (1985) *Proc. Natl. Acad. Sci.* 82: 144-148; Izant and Weintraub (1985) *Science* 229: 345-352; and Kim and Wold (1985) *Cell* 42: 129-138. Additional methods for antisense regulation are known in the art. Antisense regulation has been used to reduce or inhibit expression of plant genes in, for example in European Patent Publication No. 271988. Antisense RNA may be used to reduce gene expression to produce a visible or biochemical phenotypic change in a plant (Smith et al. (1988) *Nature* 334: 724-726; Smith et al. (1990) *Plant Mol. Biol.* 14: 369-379). In general, sense or anti-sense sequences are introduced into a cell, where they are optionally amplified, for example, by transcription. Such sequences include both simple oligonucleotide sequences and catalytic sequences such as ribozymes.

For example, a reduction or elimination of expression (i.e., a "knock-out") of a transcription factor or transcription factor homolog polypeptide in a transgenic plant, e.g., to modify a plant trait, can be obtained by introducing an antisense construct corresponding to the polypeptide of interest as a cDNA. For antisense suppression, the transcription factor or homolog cDNA is arranged in reverse orientation (with respect to the coding sequence) relative to the promoter sequence in the expression vector. The introduced sequence need not be the full length cDNA or gene, and need not be identical to the cDNA or gene found in the plant type to be transformed. Typically, the antisense sequence need only be capable of hybridizing to the target gene or RNA of interest. Thus, where the introduced sequence is of shorter length, a higher degree of homology to the endogenous transcription factor sequence will be needed for effective antisense suppression. While antisense sequences of various lengths can be utilized, preferably, the introduced antisense sequence in the vector will be at least 30 nucleotides in length, and improved antisense suppression will typically be observed as the length of the antisense sequence increases. Preferably, the length of the antisense sequence in the vector will be greater than 100 nucleotides. Transcription of an antisense construct as described results in the production of RNA molecules that are the reverse complement of mRNA molecules transcribed from the endogenous transcription factor gene in

the plant cell.

Suppression of endogenous transcription factor gene expression can also be achieved using a ribozyme. Ribozymes are RNA molecules that possess highly specific endoribonuclease activity. The production and use of ribozymes are disclosed in U.S. Patent No. 4,987,071 and U.S. Patent No. 5,543,508. Synthetic ribozyme sequences including antisense RNAs can be used to confer RNA cleaving activity on the antisense RNA, such that endogenous mRNA molecules that hybridize to the antisense RNA are cleaved, which in turn leads to an enhanced antisense inhibition of endogenous gene expression.

Vectors in which RNA encoded by a transcription factor or transcription factor homolog cDNA is over-expressed can also be used to obtain co-suppression of a corresponding endogenous gene, for example, in the manner described in U.S. Patent No. 5,231,020 to Jorgensen. Such co-suppression (also termed sense suppression) does not require that the entire transcription factor cDNA be introduced into the plant cells, nor does it require that the introduced sequence be exactly identical to the endogenous transcription factor gene of interest. However, as with antisense suppression, the suppressive efficiency will be enhanced as specificity of hybridization is increased, e.g., as the introduced sequence is lengthened, and/or as the sequence similarity between the introduced sequence and the endogenous transcription factor gene is increased.

Vectors expressing an untranslatable form of the transcription factor mRNA, e.g., sequences comprising one or more stop codon, or nonsense mutation) can also be used to suppress expression of an endogenous transcription factor, thereby reducing or eliminating its activity and modifying one or more traits. Methods for producing such constructs are described in U.S. Patent No. 5,583,021. Preferably, such constructs are made by introducing a premature stop codon into the transcription factor gene. Alternatively, a plant trait can be modified by gene silencing using double-strand RNA (Sharp (1999) *Genes and Development* 13: 139-141). Another method for abolishing the expression of a gene is by insertion mutagenesis using the T-DNA of *Agrobacterium tumefaciens*. After generating the insertion mutants, the mutants can be screened to identify those containing the insertion in a transcription factor or transcription factor homolog gene. Plants containing a single transgene insertion event at the desired gene can be crossed to generate homozygous plants for the mutation. Such methods are well known to those of skill in the art (See for example Koncz et al. (1992) *Methods in Arabidopsis Research*, World Scientific Publishing Co. Pte. Ltd., River Edge NJ).

Suppression of endogenous transcription factor gene expression can also be achieved using RNA interference, or RNAi. RNAi is a post-transcriptional, targeted gene-silencing technique that uses double-stranded RNA (dsRNA) to incite degradation of messenger RNA (mRNA) containing the same sequence as the dsRNA (Constans, (2002) *The Scientist* 16:36). Small interfering RNAs, or siRNAs are

produced in at least two steps: an endogenous ribonuclease cleaves longer dsRNA into shorter, 21-23 nucleotide-long RNAs. The siRNA segments then mediate the degradation of the target mRNA (Zamore, (2001) *Nature Struct. Biol.*, 8:746-50). RNAi has been used for gene function determination in a manner similar to antisense oligonucleotides (Constans, (2002) *The Scientist* 16:36). Expression vectors that continually express siRNAs in transiently and stably transfected have been engineered to express small hairpin RNAs (shRNAs), which get processed in vivo into siRNAs-like molecules capable of carrying out gene-specific silencing (Brummelkamp et al., (2002) *Science* 296:550-553, and Paddison, et al. (2002) *Genes & Dev.* 16:948-958). Post-transcriptional gene silencing by double-stranded RNA is discussed in further detail by Hammond et al. (2001) *Nature Rev Gen* 2: 110-119, Fire et al. (1998) *Nature* 391: 806-811 and Timmons and Fire (1998) *Nature* 395: 854.

Alternatively, a plant phenotype can be altered by eliminating an endogenous gene, such as a transcription factor or transcription factor homolog, e.g., by homologous recombination (Kempin et al. (1997) *Nature* 389: 802-803).

A plant trait can also be modified by using the Cre-lox system (for example, as described in US Pat. No. 5,658,772). A plant genome can be modified to include first and second lox sites that are then contacted with a Cre recombinase. If the lox sites are in the same orientation, the intervening DNA sequence between the two sites is excised. If the lox sites are in the opposite orientation, the intervening sequence is inverted.

The polynucleotides and polypeptides of this invention can also be expressed in a plant in the absence of an expression cassette by manipulating the activity or expression level of the endogenous gene by other means, such as, for example, by ectopically expressing a gene by T-DNA activation tagging (Ichikawa et al. (1997) *Nature* 390 698-701; Kakimoto et al. (1996) *Science* 274: 982-985). This method entails transforming a plant with a gene tag containing multiple transcriptional enhancers and once the tag has inserted into the genome, expression of a flanking gene coding sequence becomes deregulated. In another example, the transcriptional machinery in a plant can be modified so as to increase transcription levels of a polynucleotide of the invention (See, for example, PCT Publications WO 96/06166 and WO 98/53057 which describe the modification of the DNA-binding specificity of zinc finger proteins by changing particular amino acids in the DNA-binding motif).

The transgenic plant can also include the machinery necessary for expressing or altering the activity of a polypeptide encoded by an endogenous gene, for example, by altering the phosphorylation state of the polypeptide to maintain it in an activated state.

Transgenic plants (or plant cells, or plant explants, or plant tissues) incorporating the polynucleotides of the invention and/or expressing the polypeptides of the invention can be produced by a

variety of well established techniques as described above. Following construction of a vector, most typically an expression cassette, including a polynucleotide, e.g., encoding a transcription factor or transcription factor homolog, of the invention, standard techniques can be used to introduce the polynucleotide into a plant, a plant cell, a plant explant or a plant tissue of interest. Optionally, the plant cell, explant or tissue can be regenerated to produce a transgenic plant.

The plant can be any higher plant, including gymnosperms, monocotyledonous and dicotyledonous plants. Suitable protocols are available for *Leguminosae* (alfalfa, soybean, clover, etc.), *Umbelliferae* (carrot, celery, parsnip), *Cruciferae* (cabbage, radish, rapeseed, broccoli, etc.), *Curcubitaceae* (melons and cucumber), *Gramineae* (wheat, corn, rice, barley, millet, etc.), *Solanaceae* (potato, tomato, tobacco, peppers, etc.), and various other crops. See protocols described in Ammirato et al., Editors, (1984) Handbook of Plant Cell Culture – Crop Species, Macmillan Publ. Co., New York NY; Shimamoto et al. (1989) *Nature* 338: 274-276; Fromm et al. (1990) *Bio/Technol.* 8: 833-839; and Vasil et al. (1990) *Bio/Technol.* 8: 429-434.

Transformation and regeneration of both monocotyledonous and dicotyledonous plant cells are now routine, and the selection of the most appropriate transformation technique will be determined by the practitioner. The choice of method will vary with the type of plant to be transformed; those skilled in the art will recognize the suitability of particular methods for given plant types. Suitable methods can include, but are not limited to: electroporation of plant protoplasts; liposome-mediated transformation; polyethylene glycol (PEG) mediated transformation; transformation using viruses; micro-injection of plant cells; micro-projectile bombardment of plant cells; vacuum infiltration; and *Agrobacterium tumefaciens*-mediated transformation. Transformation means introducing a nucleotide sequence into a plant in a manner to cause stable or transient expression of the sequence.

Successful examples of the modification of plant characteristics by transformation with cloned sequences which serve to illustrate the current knowledge in this field of technology, and which are herein incorporated by reference, include: U.S. Patent Nos. 5,571,706; 5,677,175; 5,510,471; 5,750,386; 5,597,945; 5,589,615; 5,750,871; 5,268,526; 5,780,708; 5,538,880; 5,773,269; 5,736,369 and 5,610,042.

Following transformation, plants are preferably selected using a dominant selectable marker incorporated into the transformation vector. Typically, such a marker will confer antibiotic or herbicide resistance on the transformed plants, and selection of transformants can be accomplished by exposing the plants to appropriate concentrations of the antibiotic or herbicide.

After transformed plants are selected and grown to maturity, those plants showing a modified trait are identified. The modified trait can be any of those traits described above. Additionally, to confirm that the modified trait is due to changes in expression levels or activity of the polypeptide or polynucleotide of

the invention can be determined by analyzing mRNA expression using Northern blots, RT-PCR or microarrays, or protein expression using immunoblots or Western blots or gel shift assays.

Integrated Systems – Sequence Identity

5 Additionally, the present invention may be an integrated system, computer or computer readable medium that comprises an instruction set for determining the identity of one or more sequences in a database. In addition, the instruction set can be used to generate or identify sequences that meet any specified criteria. Furthermore, the instruction set may be used to associate or link certain functional benefits, such improved characteristics, with one or more identified sequence.

10 For example, the instruction set can include, e.g., a sequence comparison or other alignment program, e.g., an available program such as, for example, the Wisconsin Package Version 10.0, such as BLAST, FASTA, PILEUP, FINDPATTERNS or the like (GCG, Madison, WI). Public sequence databases such as GenBank, EMBL, Swiss-Prot and PIR or private sequence databases such as PHYTOSEQ sequence database (Incyte Genomics, Palo Alto CA) can be searched.

15 Alignment of sequences for comparison can be conducted by the local homology algorithm of Smith and Waterman (1981) *Adv. Appl. Math.* 2: 482-489, by the homology alignment algorithm of Needleman and Wunsch (1970) *J. Mol. Biol.* 48: 443-453, by the search for similarity method of Pearson and Lipman (1988) *Proc. Natl. Acad. Sci.* 85: 2444-2448, by computerized implementations of these algorithms. After alignment, sequence comparisons between two (or more) polynucleotides or
20 polypeptides are typically performed by comparing sequences of the two sequences over a comparison window to identify and compare local regions of sequence similarity. The comparison window can be a segment of at least about 20 contiguous positions, usually about 50 to about 200, more usually about 100 to about 150 contiguous positions. A description of the method is provided in Ausubel et al. *supra*.

 A variety of methods for determining sequence relationships can be used, including manual
25 alignment and computer assisted sequence alignment and analysis. This later approach is a preferred approach in the present invention, due to the increased throughput afforded by computer assisted methods. As noted above, a variety of computer programs for performing sequence alignment are available, or can be produced by one of skill.

 One example algorithm that is suitable for determining percent sequence identity and sequence
30 similarity is the BLAST algorithm, which is described in Altschul et al. (1990) *J. Mol. Biol.* 215: 403-410. Software for performing BLAST analyses is publicly available, e.g., through the National Library of Medicine's National Center for Biotechnology Information (ncbi.nlm.nih; see at world wide web (www) National Institutes of Health US government (gov) website). This algorithm involves first identifying high

scoring sequence pairs (HSPs) by identifying short words of length W in the query sequence, which either match or satisfy some positive-valued threshold score T when aligned with a word of the same length in a database sequence. T is referred to as the neighborhood word score threshold (Altschul et al. *supra*). These initial neighborhood word hits act as seeds for initiating searches to find longer HSPs containing them.

5 The word hits are then extended in both directions along each sequence for as far as the cumulative alignment score can be increased. Cumulative scores are calculated using, for nucleotide sequences, the parameters M (reward score for a pair of matching residues; always > 0) and N (penalty score for mismatching residues; always < 0). For amino acid sequences, a scoring matrix is used to calculate the cumulative score. Extension of the word hits in each direction are halted when: the cumulative alignment
10 score falls off by the quantity X from its maximum achieved value; the cumulative score goes to zero or below, due to the accumulation of one or more negative-scoring residue alignments; or the end of either sequence is reached. The BLAST algorithm parameters W , T , and X determine the sensitivity and speed of the alignment. The BLASTN program (for nucleotide sequences) uses as defaults a wordlength (W) of 11, an expectation (E) of 10, a cutoff of 100, $M=5$, $N=-4$, and a comparison of both strands. For amino
15 acid sequences, the BLASTP program uses as defaults a wordlength (W) of 3, an expectation (E) of 10, and the BLOSUM62 scoring matrix (*see* Henikoff and Henikoff (1992) *Proc. Natl. Acad. Sci.* 89: 10915-10919). Unless otherwise indicated, "sequence identity" here refers to the % sequence identity generated from a tblastx using the NCBI version of the algorithm at the default settings using gapped alignments with the filter "off" (*see*, for example, NIH NLM NCBI website at ncbi.nlm.nih, *supra*).

20 In addition to calculating percent sequence identity, the BLAST algorithm also performs a statistical analysis of the similarity between two sequences (*see*, for example, Karlin and Altschul (1993) *Proc. Natl. Acad. Sci.* 90: 5873-5787). One measure of similarity provided by the BLAST algorithm is the smallest sum probability ($P(N)$), which provides an indication of the probability by which a match between two nucleotide or amino acid sequences would occur by chance. For example, a nucleic acid is
25 considered similar to a reference sequence (and, therefore, in this context, homologous) if the smallest sum probability in a comparison of the test nucleic acid to the reference nucleic acid is less than about 0.1, or less than about 0.01, and or even less than about 0.001. An additional example of a useful sequence alignment algorithm is PILEUP. PILEUP creates a multiple sequence alignment from a group of related sequences using progressive, pairwise alignments. The program can align, for example, up to 300
30 sequences of a maximum length of 5,000 letters.

The integrated system, or computer typically includes a user input interface allowing a user to selectively view one or more sequence records corresponding to the one or more character strings, as well as an instruction set which aligns the one or more character strings with each other or with an additional

character string to identify one or more region of sequence similarity. The system may include a link of one or more character strings with a particular phenotype or gene function. Typically, the system includes a user readable output element that displays an alignment produced by the alignment instruction set.

The methods of this invention can be implemented in a localized or distributed computing environment. In a distributed environment, the methods may implemented on a single computer comprising multiple processors or on a multiplicity of computers. The computers can be linked, e.g. through a common bus, but more preferably the computer(s) are nodes on a network. The network can be a generalized or a dedicated local or wide-area network and, in certain preferred embodiments, the computers may be components of an intra-net or an internet.

Thus, the invention provides methods for identifying a sequence similar or homologous to one or more polynucleotides as noted herein, or one or more target polypeptides encoded by the polynucleotides, or otherwise noted herein and may include linking or associating a given plant phenotype or gene function with a sequence. In the methods, a sequence database is provided (locally or across an inter or intra net) and a query is made against the sequence database using the relevant sequences herein and associated plant phenotypes or gene functions.

Any sequence herein can be entered into the database, before or after querying the database. This provides for both expansion of the database and, if done before the querying step, for insertion of control sequences into the database. The control sequences can be detected by the query to ensure the general integrity of both the database and the query. As noted, the query can be performed using a web browser based interface. For example, the database can be a centralized public database such as those noted herein, and the querying can be done from a remote terminal or computer across an internet or intranet.

Any sequence herein can be used to identify a similar, homologous, paralogous, or orthologous sequence in another plant. This provides means for identifying endogenous sequences in other plants that may be useful to alter a trait of progeny plants, which results from crossing two plants of different strain.

For example, sequences that encode an ortholog of any of the sequences herein that naturally occur in a plant with a desired trait can be identified using the sequences disclosed herein. The plant is then crossed with a second plant of the same species but which does not have the desired trait to produce progeny which can then be used in further crossing experiments to produce the desired trait in the second plant.

Therefore the resulting progeny plant contains no transgenes; expression of the endogenous sequence may also be regulated by treatment with a particular chemical or other means, such as EMR. Some examples of such compounds well known in the art include: ethylene; cytokinins; phenolic compounds, which stimulate the transcription of the genes needed for infection; specific monosaccharides and acidic environments which potentiate vir gene induction; acidic polysaccharides which induce one or more

chromosomal genes; and opines; other mechanisms include light or dark treatment (for a review of examples of such treatments, see, Winans (1992) *Microbiol. Rev.* 56: 12-31; Eyal et al. (1992) *Plant Mol. Biol.* 19: 589-599; Chrispeels et al. (2000) *Plant Mol. Biol.* 42: 279-290; Piazza et al. (2002) *Plant Physiol.* 128: 1077-1086).

- 5 Table 5 lists sequences discovered to be orthologous to a number of representative transcription factors of the present invention. The column headings include the transcription factors listed by (a) the SEQ ID NO: of the ortholog or nucleotide encoding the ortholog; (b) the Sequence Identifier or GenBank Accession Number; (c) the species from which the orthologs to the transcription factors are derived; and (d) the smallest sum probability during by BLAST analysis.

10

Table 5. Paralogs and Orthologs and Other Related Genes of Representative *Arabidopsis* Transcription Factor Genes identified using BLAST

SEQ ID NO: of Ortholog or Nucleotide Encoding Ortholog	GID No.	Sequence Identifier or Accession Number	Species from Which Ortholog is Derived	Smallest Sum Probability to <i>Arabidopsis</i> Polynucleotide Sequence
3	G1067		<i>Arabidopsis thaliana</i>	
5	G2153		<i>Arabidopsis thaliana</i>	
7	G2156		<i>Arabidopsis thaliana</i>	
41	G1069		<i>Arabidopsis thaliana</i>	5e-90**
43	G1945		<i>Arabidopsis thaliana</i>	5e-51**
45	G2155		<i>Arabidopsis thaliana</i>	6e-43**
47	G1070		<i>Arabidopsis thaliana</i>	5e-70**
49	G2657		<i>Arabidopsis thaliana</i>	3e-70†
51	G1075		<i>Arabidopsis thaliana</i>	8e-72**
53	G1076		<i>Arabidopsis thaliana</i>	9e-74**
9	G3399	AP004165	<i>Oryza sativa</i> (japonica cultivar-group)	1e-81†
11	G3407	AP004635	<i>Oryza sativa</i>	5e-90†
13	G3456	BM525692	<i>Glycine max</i>	2e-87**
39	G3556		<i>Oryza sativa</i>	7e-67††
15	G3459	C33095_1	<i>Glycine max</i>	6e-67††
17	G3460	C33095_2	<i>Glycine max</i>	1e-66*
65		BH566718	<i>Brassica oleracea</i>	1e-129**
67		BH685875	<i>Brassica oleracea</i>	1e-124†
		BZ432677	<i>Brassica oleracea</i>	1e-113**
		BZ433664	<i>Brassica oleracea</i>	1e-107†
		BH730050	<i>Brassica oleracea</i>	1e-104†

		AP004971	<i>Lotus corniculatus</i> var. <i>japonicus</i>	3e-91**
		CC729476	<i>Zea mays</i>	1e-83**
21	G3403	AP004020	<i>Oryza sativa</i> (<i>japonica</i> cultivar-group)	2e-81**
		AAAA01000486	<i>Oryza sativa</i> (<i>indica</i> cultivar-group)	7e-80*
		CB003423	<i>Vitis vinifera</i>	2e-76*
		CC645378	<i>Zea mays</i>	4e-75*
23	G3458	C32394_2	<i>Glycine max</i>	9e-73**
25	G3406	AL662981	<i>Oryza sativa</i>	7e-73*
		BQ785950	<i>Glycine max</i>	3e-73*
		BH975957	<i>Brassica oleracea</i>	9e-72*
		BQ865858	<i>Lactuca sativa</i>	7e-72*
		CB891166	<i>Medicago truncatula</i>	5e-72*
		CF229888	<i>Populus x canescens</i>	2e-71*
		BQ863249	<i>Lactuca sativa</i>	2e-71*
		BG134451	<i>Lycopersicon esculentum</i>	3e-70*
27	G3405	AP005653	<i>Oryza sativa</i> (<i>japonica</i> cultivar-group)	1e-69**
29	G3400	AP005477	<i>Oryza sativa</i> (<i>japonica</i> cultivar-group)	2e-67*
31	G3404	AP003526	<i>Oryza sativa</i> (<i>japonica</i> cultivar-group)	2e-67*
		AP004971	<i>Lotus corniculatus</i> var. <i>japonicus</i>	7e-66*
		BM110212	<i>Solanum tuberosum</i>	8e-65*
33	G3407	AP004635	<i>Oryza sativa</i> (<i>japonica</i> cultivar-group)	6e-63*
		AC124953	<i>Medicago truncatula</i>	2e-63*
35	G3462	BI321563	<i>Glycine max</i>	3e-61*
		BH660108	<i>Brassica oleracea</i>	2e-61†
		BQ838600	<i>Triticum aestivum</i>	2e-59*
		CD825510	<i>Brassica napus</i>	7e-58†
		BF254863	<i>Hordeum vulgare</i>	1e-56*

19	G3408	AP005755	<i>Oryza sativa</i>	5e-43††
37	G3401	AAAA01017331 SC17331 AP004587	<i>Oryza sativa (japonica)</i> cultivar-group	9e-42*

* Smallest sum probability comparison to G1073

† Smallest sum probability comparison to G1067

** Smallest sum probability comparison to G2153

†† Smallest sum probability comparison to 2156

5

Molecular Modeling

Another means that may be used to confirm the utility and function of transcription factor sequences that are orthologous or paralogous to presently disclosed transcription factors is through the use of molecular modeling software. Molecular modeling is routinely used to predict polypeptide structure, and a variety of protein structure modeling programs, such as "Insight II" (Accelrys, Inc.) are commercially available for this purpose. Modeling can thus be used to predict which residues of a polypeptide can be changed without altering function (Crameri et al. (2003) U.S. Patent No. 6, 521, 453). Thus, polypeptides that are sequentially similar can be shown to have a high likelihood of similar function by their structural similarity, which may, for example, be established by comparison of regions of superstructure. The relative tendencies of amino acids to form regions of superstructure (for example, helixes and α -sheets) are well established. For example, O'Neil et al. ((1990) *Science* 250: 646-651) have discussed in detail the helix forming tendencies of amino acids. Tables of relative structure forming activity for amino acids can be used as substitution tables to predict which residues can be functionally substituted in a given region, for example, in DNA-binding domains of known transcription factors and

10

15

20

equivalogs. Homologs that are likely to be functionally similar can then be identified.

Of particular interest is the structure of a transcription factor in the region of its conserved domains, such as those identified in Table 1. Structural analyses may be performed by comparing the structure of the known transcription factor around its conserved domain with those of orthologs and paralogs. Analysis of a number of polypeptides within a transcription factor group or clade, including the functionally or sequentially similar polypeptides provided in the Sequence Listing, may also provide an understanding of structural elements required to regulate transcription within a given family.

25

EXAMPLES

It is to be understood that this invention is not limited to the particular devices, machines, materials and methods described. Although particular embodiments are described, equivalent embodiments may be used to practice the invention. The described embodiments are not intended to limit

30

the scope of the invention, which is limited only by the appended claims. The examples below are provided to enable the subject invention and are not included for the purpose of limiting the invention.

The invention, now being generally described, will be more readily understood by reference to the following examples, which are included merely for purposes of illustration of certain aspects and embodiments of the present invention and are not intended to limit the invention. It will be recognized by one of skill in the art that a transcription factor that is associated with a particular first trait may also be associated with at least one other, unrelated and inherent second trait which was not predicted by the first trait.

Example I: Full Length Gene Identification and Cloning

Putative transcription factor sequences (genomic or ESTs) related to known transcription factors were identified in the *Arabidopsis thaliana* GenBank database using the tblastn sequence analysis program using default parameters and a P-value cutoff threshold of -4 or -5 or lower, depending on the length of the query sequence. Putative transcription factor sequence hits were then screened to identify those containing particular sequence strings. If the sequence hits contained such sequence strings, the sequences were confirmed as transcription factors.

Alternatively, *Arabidopsis thaliana* cDNA libraries derived from different tissues or treatments, or genomic libraries were screened to identify novel members of a transcription family using a low stringency hybridization approach. Probes were synthesized using gene specific primers in a standard PCR reaction (annealing temperature 60°C) and labeled with ^{32}P dCTP using the High Prime DNA Labeling Kit (Boehringer Mannheim Corp. (now Roche Diagnostics Corp., Indianapolis, IN). Purified radiolabelled probes were added to filters immersed in Church hybridization medium (0.5 M NaPO_4 pH 7.0, 7% SDS, 1% w/v bovine serum albumin) and hybridized overnight at 60°C with shaking. Filters were washed two times for 45 to 60 minutes with 1xSSC, 1% SDS at 60°C .

To identify additional sequence 5' or 3' of a partial cDNA sequence in a cDNA library, 5' and 3' rapid amplification of cDNA ends (RACE) was performed using the MARATHON cDNA amplification kit (Clontech, Palo Alto, CA). Generally, the method entailed first isolating poly(A) mRNA, performing first and second strand cDNA synthesis to generate double stranded cDNA, blunting cDNA ends, followed by ligation of the MARATHON Adaptor to the cDNA to form a library of adaptor-ligated ds cDNA.

Gene-specific primers were designed to be used along with adaptor specific primers for both 5' and 3' RACE reactions. Nested primers, rather than single primers, were used to increase PCR specificity. Using 5' and 3' RACE reactions, 5' and 3' RACE fragments were obtained, sequenced and cloned. The process can be repeated until 5' and 3' ends of the full-length gene were identified. Then the full-length

cDNA was generated by PCR using primers specific to 5' and 3' ends of the gene by end-to-end PCR.

Example II: Construction of Expression Vectors

The sequence was amplified from a genomic or cDNA library using primers specific to sequences upstream and downstream of the coding region. The expression vector was pMEN20 or pMEN65, which are both derived from pMON316 (Sanders et al. (1987) *Nucleic Acids Res.* 15:1543-1558) and contain the CaMV 35S promoter to express transgenes. To clone the sequence into the vector, both pMEN20 and the amplified DNA fragment were digested separately with SalI and NotI restriction enzymes at 37° C for 2 hours. The digestion products were subject to electrophoresis in a 0.8% agarose gel and visualized by ethidium bromide staining. The DNA fragments containing the sequence and the linearized plasmid were excised and purified by using a QIAQUICK gel extraction kit (Qiagen, Valencia, CA). The fragments of interest were ligated at a ratio of 3:1 (vector to insert). Ligation reactions using T4 DNA ligase (New England Biolabs, Beverly MA) were carried out at 16° C for 16 hours. The ligated DNAs were transformed into competent cells of the *E. coli* strain DH5alpha by using the heat shock method. The transformations were plated on LB plates containing 50 mg/l kanamycin (Sigma Chemical Co. St. Louis MO). Individual colonies were grown overnight in five milliliters of LB broth containing 50 mg/l kanamycin at 37° C. Plasmid DNA was purified by using Qiaquick Mini Prep kits (Qiagen, Valencia CA).

Example III: Transformation of *Agrobacterium* with the Expression Vector

After the plasmid vector containing the gene was constructed, the vector was used to transform *Agrobacterium tumefaciens* cells expressing the gene products. The stock of *Agrobacterium tumefaciens* cells for transformation were made as described by Nagel et al. (1990) *FEMS Microbiol Letts.* 67: 325-328. *Agrobacterium* strain ABI was grown in 250 ml LB medium (Sigma) overnight at 28° C with shaking until an absorbance over 1 cm at 600 nm (A_{600}) of 0.5 – 1.0 was reached. Cells were harvested by centrifugation at 4,000 x g for 15 min at 4° C. Cells were then resuspended in 250 µl chilled buffer (1 mM HEPES, pH adjusted to 7.0 with KOH). Cells were centrifuged again as described above and resuspended in 125 µl chilled buffer. Cells were then centrifuged and resuspended two more times in the same HEPES buffer as described above at a volume of 100 µl and 750 µl, respectively. Resuspended cells were then distributed into 40 µl aliquots, quickly frozen in liquid nitrogen, and stored at -80° C.

Agrobacterium cells were transformed with plasmids prepared as described above following the protocol described by Nagel et al. (*supra*). For each DNA construct to be transformed, 50 – 100 ng DNA (generally resuspended in 10 mM Tris-HCl, 1 mM EDTA, pH 8.0) was mixed with 40 µl of *Agrobacterium* cells. The DNA/cell mixture was then transferred to a chilled cuvette with a 2mm

electrode gap and subject to a 2.5 kV charge dissipated at 25 μ F and 200 μ F using a Gene Pulser II apparatus (Bio-Rad, Hercules, CA). After electroporation, cells were immediately resuspended in 1.0 ml LB and allowed to recover without antibiotic selection for 2 – 4 hours at 28° C in a shaking incubator. After recovery, cells were plated onto selective medium of LB broth containing 100 μ g/ml spectinomycin (Sigma) and incubated for 24-48 hours at 28° C. Single colonies were then picked and inoculated in fresh medium. The presence of the plasmid construct was verified by PCR amplification and sequence analysis.

Example IV: Transformation of *Arabidopsis* Plants with *Agrobacterium tumefaciens* with

Expression Vector

After transformation of *Agrobacterium tumefaciens* with plasmid vectors containing the gene, single *Agrobacterium* colonies were identified, propagated, and used to transform *Arabidopsis* plants. Briefly, 500 ml cultures of LB medium containing 50 mg/l kanamycin were inoculated with the colonies and grown at 28° C with shaking for 2 days until an optical absorbance at 600 nm wavelength over 1 cm (A_{600}) of > 2.0 is reached. Cells were then harvested by centrifugation at 4,000 x g for 10 min, and resuspended in infiltration medium (1/2 X Murashige and Skoog salts (Sigma), 1 X Gamborg's B-5 vitamins (Sigma), 5.0% (w/v) sucrose (Sigma), 0.044 μ M benzylamino purine (Sigma), 200 μ l/l Silwet L-77 (Lehle Seeds) until an A_{600} of 0.8 was reached.

Prior to transformation, *Arabidopsis thaliana* seeds (ecotype Columbia) were sown at a density of ~10 plants per 4" pot onto Pro-Mix BX potting medium (Hummert International) covered with fiberglass mesh (18 mm X 16 mm). Plants were grown under continuous illumination (50-75 μ E/m²/sec) at 22-23° C with 65-70% relative humidity. After about 4 weeks, primary inflorescence stems (bolts) are cut off to encourage growth of multiple secondary bolts. After flowering of the mature secondary bolts, plants were prepared for transformation by removal of all siliques and opened flowers.

The pots were then immersed upside down in the mixture of *Agrobacterium* infiltration medium as described above for 30 sec, and placed on their sides to allow draining into a 1' x 2' flat surface covered with plastic wrap. After 24 h, the plastic wrap was removed and pots are turned upright. The immersion procedure was repeated one week later, for a total of two immersions per pot. Seeds were then collected from each transformation pot and analyzed following the protocol described below.

Example V: Identification of *Arabidopsis* Primary Transformants

Seeds collected from the transformation pots were sterilized essentially as follows. Seeds were dispersed into in a solution containing 0.1% (v/v) Triton X-100 (Sigma) and sterile water and washed by

shaking the suspension for 20 min. The wash solution was then drained and replaced with fresh wash solution to wash the seeds for 20 min with shaking. After removal of the ethanol/detergent solution, a solution containing 0.1% (v/v) Triton X-100 and 30% (v/v) bleach (CLOROX; Clorox Corp. Oakland CA) was added to the seeds, and the suspension was shaken for 10 min. After removal of the bleach/detergent solution, seeds were then washed five times in sterile distilled water. The seeds were stored in the last wash water at 4° C for 2 days in the dark before being plated onto antibiotic selection medium (1 X Murashige and Skoog salts (pH adjusted to 5.7 with 1M KOH), 1 X Gamborg's B-5 vitamins, 0.9% phytagar (Life Technologies), and 50 mg/l kanamycin). Seeds were germinated under continuous illumination (50-75 $\mu\text{E}/\text{m}^2/\text{sec}$) at 22-23° C. After 7-10 days of growth under these conditions, kanamycin resistant primary transformants (T_1 generation) were visible and obtained. These seedlings were transferred first to fresh selection plates where the seedlings continued to grow for 3-5 more days, and then to soil (Pro-Mix BX potting medium).

Primary transformants were crossed and progeny seeds (T_2) collected; kanamycin resistant seedlings were selected and analyzed. The expression levels of the recombinant polynucleotides in the transformants varies from about a 5% expression level increase to a least a 100% expression level increase. Similar observations are made with respect to polypeptide level expression.

Example VI: Identification of *Arabidopsis* Plants with Transcription Factor Gene Knockouts

The screening of insertion mutagenized *Arabidopsis* collections for null mutants in a known target gene was essentially as described in Krysan et al. (1999) *Plant Cell* 11: 2283-2290. Briefly, gene-specific primers, nested by 5-250 base pairs to each other, were designed from the 5' and 3' regions of a known target gene. Similarly, nested sets of primers were also created specific to each of the T-DNA or transposon ends (the "right" and "left" borders). All possible combinations of gene specific and T-DNA/transposon primers were used to detect by PCR an insertion event within or close to the target gene. The amplified DNA fragments were then sequenced which allows the precise determination of the T-DNA/transposon insertion point relative to the target gene. Insertion events within the coding or intervening sequence of the genes were deconvoluted from a pool comprising a plurality of insertion events to a single unique mutant plant for functional characterization. The method is described in more detail in Yu and Adam, US Application Serial No. 09/177,733 filed October 23, 1998.

Example VII: Identification of Modified Phenotypes in Overexpressing or Knockout Plants

Experiments were performed to identify those transformants or knockouts that exhibited modified biochemical characteristics. Among the biochemicals that were assayed were insoluble sugars, such as arabinose, fucose, galactose, mannose, rhamnose or xylose or the like; prenyl lipids, such as lutein, beta-

carotene, xanthophyll-1, xanthophyll-2, chlorophylls A or B, or alpha-, delta- or gamma-tocopherol or the like; fatty acids, such as 16:0 (palmitic acid), 16:1 (palmitoleic acid), 18:0 (stearic acid), 18:1 (oleic acid), 18:2 (linoleic acid), 20:0, 18:3 (linolenic acid), 20:1 (eicosenoic acid), 20:2, 22:1 (erucic acid) or the like; waxes, such as by altering the levels of C29, C31, or C33 alkanes; sterols, such as brassicasterol, campesterol, stigmasterol, sitosterol or stigmasterol or the like, glucosinolates, protein or oil levels.

Fatty acids were measured using two methods depending on whether the tissue was from leaves or seeds. For leaves, lipids were extracted and esterified with hot methanolic H_2SO_4 and partitioned into hexane from methanolic brine. For seed fatty acids, seeds were pulverized and extracted in methanol:heptane:toluene:2,2-dimethoxypropane: H_2SO_4 (39:34:20:5:2) for 90 minutes at 80°C. After cooling to room temperature the upper phase, containing the seed fatty acid esters, was subjected to GC analysis. Fatty acid esters from both seed and leaf tissues were analyzed with a SUPELCO SP-2330 column (Supelco, Bellefonte, PA).

Glucosinolates were purified from seeds or leaves by first heating the tissue at 95°C for 10 minutes. Preheated ethanol:water (50:50) is added and after heating at 95°C for a further 10 minutes, the extraction solvent is applied to a DEAE Sephadex column (Pharmacia) which had been previously equilibrated with 0.5 M pyridine acetate. Desulfoglucosinolates were eluted with 300 ul water and analyzed by reverse phase HPLC monitoring at 226 nm.

For wax alkanes, samples were extracted using an identical method as fatty acids and extracts were analyzed on a HP 5890 GC coupled with a 5973 MSD. Samples were chromatographically isolated on a J&W DB35 mass spectrometer (J&W Scientific Agilent Technologies, Folsom, CA).

To measure prenyl lipid levels, seeds or leaves were pulverized with 1 to 2% pyrogallol as an antioxidant. For seeds, extracted samples were filtered and a portion removed for tocopherol and carotenoid/chlorophyll analysis by HPLC. The remaining material was saponified for sterol determination. For leaves, an aliquot was removed and diluted with methanol and chlorophyll A, chlorophyll B, and total carotenoids measured by spectrophotometry by determining optical absorbance at 665.2 nm, 652.5 nm, and 470 nm. An aliquot was removed for tocopherol and carotenoid/chlorophyll composition by HPLC using a Waters μ Bondapak C18 column (4.6 mm x 150 mm). The remaining methanolic solution was saponified with 10% KOH at 80°C for one hour. The samples were cooled and diluted with a mixture of methanol and water. A solution of 2% methylene chloride in hexane was mixed in and the samples were centrifuged. The aqueous methanol phase was again re-extracted 2% methylene chloride in hexane and, after centrifugation, the two upper phases were combined and evaporated. 2% methylene chloride in hexane was added to the tubes and the samples were then extracted with one ml of water. The upper phase

was removed, dried, and resuspended in 400 μ l of 2% methylene chloride in hexane and analyzed by gas chromatography using a 50 m DB-5ms (0.25 mm ID, 0.25 μ m phase, J&W Scientific).

Insoluble sugar levels were measured by the method essentially described by Reiter et al. (1999), *Plant J.* 12: 335-345. This method analyzes the neutral sugar composition of cell wall polymers found in *Arabidopsis* leaves. Soluble sugars were separated from sugar polymers by extracting leaves with hot 70% ethanol. The remaining residue containing the insoluble polysaccharides was then acid hydrolyzed with allose added as an internal standard. Sugar monomers generated by the hydrolysis were then reduced to the corresponding alditols by treatment with NaBH₄, then were acetylated to generate the volatile alditol acetates which were then analyzed by GC-FID. Identity of the peaks was determined by comparing the retention times of known sugars converted to the corresponding alditol acetates with the retention times of peaks from wild-type plant extracts. Alditol acetates were analyzed on a Supelco SP-2330 capillary column (30 m x 250 μ m x 0.2 μ m) using a temperature program beginning at 180° C for 2 minutes followed by an increase to 220° C in 4 minutes. After holding at 220° C for 10 minutes, the oven temperature is increased to 240° C in 2 minutes and held at this temperature for 10 minutes and brought back to room temperature.

To identify plants with alterations in total seed oil or protein content, 150mg of seeds from T2 progeny plants were subjected to analysis by Near Infrared Reflectance Spectroscopy (NIRS) using a Foss NirSystems Model 6500 with a spinning cup transport system. NIRS is a non-destructive analytical method used to determine seed oil and protein composition. Infrared is the region of the electromagnetic spectrum located after the visible region in the direction of longer wavelengths. 'Near infrared' owns its name for being the infrared region near to the visible region of the electromagnetic spectrum. For practical purposes, near infrared comprises wavelengths between 800 and 2500 nm. NIRS is applied to organic compounds rich in O-H bonds (such as moisture, carbohydrates, and fats), C-H bonds (such as organic compounds and petroleum derivatives), and N-H bonds (such as proteins and amino acids). The NIRS analytical instruments operate by statistically correlating NIRS signals at several wavelengths with the characteristic or property intended to be measured. All biological substances contain thousands of C-H, O-H, and N-H bonds. Therefore, the exposure to near infrared radiation of a biological sample, such as a seed, results in a complex spectrum which contains qualitative and quantitative information about the physical and chemical composition of that sample.

The numerical value of a specific analyte in the sample, such as protein content or oil content, is mediated by a calibration approach known as chemometrics. Chemometrics applies statistical methods such as multiple linear regression (MLR), partial least squares (PLS), and principle component analysis (PCA) to the spectral data and correlates them with a physical property or other factor, that property or

factor is directly determined rather than the analyte concentration itself. The method first provides "wet chemistry" data of the samples required to develop the calibration.

Calibration of NIRS response was performed using data obtained by wet chemical analysis of a population of *Arabidopsis* ecotypes that were expected to represent diversity of oil and protein levels.

5 The exact oil composition of each ecotype used in the calibration experiment was performed using gravimetric analysis of oils extracted from seed samples (0.5 g or 1.0 g) by the accelerated solvent extraction method (ASE; Dionex Corp, Sunnyvale, CA). The extraction method was validated against certified canola samples (Community Bureau of Reference, Belgium). Seed samples from each ecotype (0.5 g or 1g) were subjected to accelerated solvent extraction and the resulting extracted oil weights
10 compared to the weight of oil recovered from canola seed that has been certified for oil content (Community Bureau of Reference). The oil calibration equation was based on 57 samples with a range of oil contents from 27.0 % to 50.8 %. To check the validity of the calibration curve, an additional set of samples was extracted by ASE and predicted using the oil calibration equation. This validation set counted 46 samples, ranging from 27.9 % to 47.5 % oil, and had a predicted standard error of performance of 0.63
15 %. The wet chemical method for protein was elemental analysis (%N X 6.0) using the average of 3 representative samples of 5 mg each validated against certified ground corn (NIST). The instrumentation was an Elementar Vario-EL III elemental analyzer operated in CNS operating mode (Elementar Analysensysteme GmbH, Hanau, Germany).

20 The protein calibration equation was based on a library of 63 samples with a range of protein contents from 17.4 % to 31.2 %. An additional set of samples was analyzed for protein by elemental analysis ($n = 57$) and scanned by NIRS in order to validate the protein prediction equation. The protein range of the validation set was from 16.8 % to 31.2 % and the standard error of prediction was 0.468 %.

NIRS analysis of *Arabidopsis* seed was carried out on between 40-300 mg experimental sample. The oil and protein contents were predicted using the respective calibration equations.

25 Data obtained from NIRS analysis was analyzed statistically using a nearest-neighbor (N-N) analysis. The N-N analysis allows removal of within-block spatial variability in a fairly flexible fashion, which does not require prior knowledge of the pattern of variability in the chamber. Ideally, all hybrids are grown under identical experimental conditions within a block (rep). In reality, even in many block designs, significant within-block variability exists. Nearest-neighbor procedures are based on assumption that
30 environmental effect of a plot is closely related to that of its neighbors. Nearest-neighbor methods use information from adjacent plots to adjust for within-block heterogeneity and so provide more precise estimates of treatment means and differences. If there is within-plot heterogeneity on a spatial scale that is larger than a single plot and smaller than the entire block, then yields from adjacent plots will be positively

correlated. Information from neighboring plots can be used to reduce or remove the unwanted effect of the spatial heterogeneity, and hence improve the estimate of the treatment effect. Data from neighboring plots can also be used to reduce the influence of competition between adjacent plots. The Papadakis N-N analysis can be used with designs to remove within-block variability that would not be removed with the standard split plot analysis (Papadakis (1973) *Inst. d'Amelior. Plantes Thessaloniki (Greece) Bull. Scientif.* No. 23; Papadakis (1984) *Proc. Acad. Athens* 59: 326-342).

Experiments were performed to identify those transformants or knockouts that exhibited modified sugar-sensing. For such studies, seeds from transformants were germinated on media containing 5% glucose or 9.4% sucrose which normally partially restrict hypocotyl elongation. Plants with altered sugar sensing may have either longer or shorter hypocotyls than normal plants when grown on this media. Additionally, other plant traits may be varied such as root mass.

Experiments may be performed to identify those transformants or knockouts that exhibited an improved pathogen tolerance. For such studies, the transformants are exposed to biotrophic fungal pathogens, such as *Erysiphe orontii*, and necrotrophic fungal pathogens, such as *Fusarium oxysporum*. *Fusarium oxysporum* isolates cause vascular wilts and damping off of various annual vegetables, perennials and weeds (Mauch-Mani and Slusarenko (1994) *Molec Plant-Microbe Interact.* 7: 378-383). For *Fusarium oxysporum* experiments, plants are grown on Petri dishes and sprayed with a fresh spore suspension of *F. oxysporum*. The spore suspension is prepared as follows: A plug of fungal hyphae from a plate culture is placed on a fresh potato dextrose agar plate and allowed to spread for one week. Five ml sterile water is then added to the plate, swirled, and pipetted into 50 ml Armstrong *Fusarium* medium. Spores are grown overnight in *Fusarium* medium and then sprayed onto plants using a Preval paint sprayer. Plant tissue is harvested and frozen in liquid nitrogen 48 hours post-infection.

Erysiphe orontii is a causal agent of powdery mildew. For *Erysiphe orontii* experiments, plants are grown approximately 4 weeks in a greenhouse under 12 hour light (20°C, ~30% relative humidity (rh)). Individual leaves are infected with *E. orontii* spores from infected plants using a camel's hair brush, and the plants are transferred to a Percival growth chamber (20°C, 80% rh.). Plant tissue is harvested and frozen in liquid nitrogen 7 days post-infection.

Botrytis cinerea is a necrotrophic pathogen. *Botrytis cinerea* is grown on potato dextrose agar under 12 hour light (20°C, ~30% relative humidity (rh)). A spore culture is made by spreading 10 ml of sterile water on the fungus plate, swirling and transferring spores to 10 ml of sterile water. The spore inoculum (approx. 105 spores/ml) is then used to spray 10 day-old seedlings grown under sterile conditions on MS (minus sucrose) media. Symptoms are evaluated every day up to approximately 1 week.

Sclerotinia sclerotiorum hyphal cultures are grown in potato dextrose broth. One gram of hyphae is ground, filtered, spun down and resuspended in sterile water. A 1:10 dilution is used to spray 10 day-old seedlings grown aseptically under a 12 hour light/dark regime on MS (minus sucrose) media. Symptoms are evaluated every day up to approximately 1 week.

Pseudomonas syringae pv *maculicola* (Psm) strain 4326 and pv *maculicola* strain 4326 was inoculated by hand at two doses. Two inoculation doses allows the differentiation between plants with enhanced susceptibility and plants with enhanced resistance to the pathogen. Plants are grown for 3 weeks in the greenhouse, then transferred to the growth chamber for the remainder of their growth. Psm ES4326 may be hand inoculated with 1 ml syringe on 3 fully-expanded leaves per plant (4 1/2 wk old), using at least 9 plants per overexpressing line at two inoculation doses, OD=0.005 and OD=0.0005. Disease scoring is performed at day 3 post-inoculation with pictures of the plants and leaves taken in parallel.

In some instances, expression patterns of the pathogen-induced genes (such as defense genes) may be monitored by microarray experiments. In these experiments, cDNAs are generated by PCR and resuspended at a final concentration of ~ 100 ng/ µl in 3X SSC or 150mM Na-phosphate (Eisen and Brown (1999) *Methods Enzymol.* 303: 179-205). The cDNAs are spotted on microscope glass slides coated with polylysine. The prepared cDNAs are aliquoted into 384 well plates and spotted on the slides using, for example, an x-y-z gantry (OmniGrid) which may be purchased from GeneMachines (Menlo Park, CA) outfitted with quill type pins which may be purchased from Telechem International (Sunnyvale, CA). After spotting, the arrays are cured for a minimum of one week at room temperature, rehydrated and blocked following the protocol recommended by Eisen and Brown (1999; *supra*).

Sample total RNA (10 µg) samples are labeled using fluorescent Cy3 and Cy5 dyes. Labeled samples are resuspended in 4X SSC/0.03% SDS/4 µg salmon sperm DNA/2 µg tRNA/ 50mM Na-pyrophosphate, heated for 95°C for 2.5 minutes, spun down and placed on the array. The array is then covered with a glass coverslip and placed in a sealed chamber. The chamber is then kept in a water bath at 62°C overnight. The arrays are washed as described in Eisen and Brown (1999, *supra*) and scanned on a General Scanning 3000 laser scanner. The resulting files are subsequently quantified using IMAGE, software (BioDiscovery, Los Angeles CA).

RT-PCR experiments may be performed to identify those genes induced after exposure to biotrophic fungal pathogens, such as *Erysiphe orontii*, necrotrophic fungal pathogens, such as *Fusarium oxysporum*, bacteria, viruses and salicylic acid, the latter being involved in a nonspecific resistance response in *Arabidopsis thaliana*. Generally, the gene expression patterns from ground plant leaf tissue is examined.

Reverse transcriptase PCR was conducted using gene specific primers within the coding region for each sequence identified. The primers were designed near the 3' region of each DNA binding sequence initially identified.

Total RNA from these ground leaf tissues was isolated using the CTAB extraction protocol. Once extracted total RNA was normalized in concentration across all the tissue types to ensure that the PCR reaction for each tissue received the same amount of cDNA template using the 28S band as reference. Poly(A+) RNA was purified using a modified protocol from the Qiagen OLIGOTEX purification kit batch protocol. cDNA was synthesized using standard protocols. After the first strand cDNA synthesis, primers for Actin 2 were used to normalize the concentration of cDNA across the tissue types. Actin 2 is found to be constitutively expressed in fairly equal levels across the tissue types being investigated.

For RT PCR, cDNA template was mixed with corresponding primers and Taq DNA polymerase. Each reaction consisted of 0.2 µl cDNA template, 2 µl 10X Tricine buffer, 2 µl 10X Tricine buffer and 16.8 µl water, 0.05 µl Primer 1, 0.05 µl, Primer 2, 0.3 µl Taq DNA polymerase and 8.6 µl water.

The 96 well plate is covered with microfilm and set in the thermocycler to start the reaction cycle.

By way of illustration, the reaction cycle may comprise the following steps:

Step 1: 93° C for 3 min;

Step 2: 93° C for 30 sec;

Step 3: 65° C for 1 min;

Step 4: 72° C for 2 min;

Steps 2, 3 and 4 are repeated for 28 cycles;

Step 5: 72° C for 5 min; and

Step 6 4° C.

To amplify more products, for example, to identify genes that have very low expression, additional steps may be performed: The following method illustrates a method that may be used in this regard. The PCR plate is placed back in the thermocycler for 8 more cycles of steps 2-4.

Step 2 93° C for 30 sec;

Step 3 65° C for 1 min;

Step 4 72° C for 2 min, repeated for 8 cycles; and

Step 5 4° C.

Eight microliters of PCR product and 1.5 µl of loading dye are loaded on a 1.2% agarose gel for analysis after 28 cycles and 36 cycles. Expression levels of specific transcripts are considered low if they were only detectable after 36 cycles of PCR. Expression levels are considered medium or high depending on the levels of transcript compared with observed transcript levels for an internal control such as actin2.

Transcript levels are determined in repeat experiments and compared to transcript levels in control (e.g., non-transformed) plants.

Modified phenotypes observed for particular overexpressor or knockout plants may include increased biomass, and/or increased or decreased abiotic stress tolerance or resistance. For a particular overexpressor that shows a less beneficial characteristic, such as reduced disease resistance or tolerance, it may be more useful to select a plant with a decreased expression of the particular transcription factor. For a particular knockout that shows a less beneficial characteristic, such as decreased abiotic stress tolerance, it may be more useful to select a plant with an increased expression of the particular transcription factor.

The transcription factor sequences of the Sequence Listing, or those in the present Tables or Figures, and their equivalents, can be used to prepare transgenic plants and plants with altered traits. The specific transgenic plants listed below are produced from the sequences of the Sequence Listing, as noted. The Sequence Listing and Table 5 provide exemplary polynucleotide and polypeptide sequences of the invention.

Example VIII: Genes that Confer Significant Improvements to Plants

Examples of genes and homologs that confer significant improvements to knockout or overexpressing plants are noted below. Experimental observations made by us with regard to specific genes whose expression has been modified in overexpressing or knock-out plants, and potential applications based on these observations, are also presented.

This example provides experimental evidence for increased biomass and abiotic stress tolerance controlled by the transcription factor polypeptides and polypeptides of the invention.

Salt stress assays are intended to find genes that confer better germination, seedling vigor or growth in high salt. Evaporation from the soil surface causes upward water movement and salt accumulation in the upper soil layer where the seeds are placed. Thus, germination normally takes place at a salt concentration much higher than the mean salt concentration of in the whole soil profile. Plants differ in their tolerance to NaCl depending on their stage of development, therefore seed germination, seedling vigor, and plant growth responses are evaluated.

Osmotic stress assays (including NaCl and mannitol assays) are intended to determine if an osmotic stress phenotype is NaCl-specific or if it is a general osmotic stress related phenotype. Plants tolerant to osmotic stress could also have more tolerance to drought and/or freezing.

Drought assays are intended to find genes that mediate better plant survival after short-term, severe water deprivation. Ion leakage will be measured if needed. Osmotic stress tolerance would also support a drought tolerant phenotype.

Temperature stress assays are intended to find genes that confer better germination, seedling vigor or plant growth under temperature stress (cold, freezing and heat).

Sugar sensing assays are intended to find genes involved in sugar sensing by germinating seeds on high concentrations of sucrose and glucose and looking for degrees of hypocotyl elongation. The

5 germination assay on mannitol controls for responses related to osmotic stress. Sugars are key regulatory molecules that affect diverse processes in higher plants including germination, growth, flowering, senescence, sugar metabolism and photosynthesis. Sucrose is the major transport form of photosynthate and its flux through cells has been shown to affect gene expression and alter storage compound accumulation in seeds (source-sink relationships). Glucose-specific hexose-sensing has also been
10 described in plants and is implicated in cell division and repression of "famine" genes (photosynthetic or glyoxylate cycles).

Germination assays followed modifications of the same basic protocol. Sterile seeds were sown on the conditional media listed below. Plates were incubated at 22° C under 24-hour light (120-130 $\mu\text{Ein}/\text{m}^2/\text{s}$) in a growth chamber. Evaluation of germination and seedling vigor was conducted 3 to 15 days
15 after planting. The basal media was 80% Murashige-Skoog medium (MS) + vitamins.

For salt and osmotic stress germination experiments, the medium was supplemented with 150 mM NaCl or 300 mM mannitol. Growth regulator sensitivity assays were performed in MS media, vitamins, and either 0.3 μM ABA, 9.4% sucrose, or 5% glucose.

20 Temperature stress cold germination experiments were carried out at 8 °C. Heat stress germination experiments were conducted at 32 °C to 37° C for 6 hours of exposure.

For stress experiments conducted with more mature plants, seeds were germinated and grown for seven days on MS + vitamins + 1% sucrose at 22 °C and then transferred to chilling and heat stress conditions. The plants were either exposed to chilling stress (6 hour exposure to 4-8° C), or heat stress (32° C was applied for five days, after which the plants were transferred back 22 °C for recovery and
25 evaluated after 5 days relative to controls not exposed to the depressed or elevated temperature).

Results:

G1073 (SEQ ID NOs: 1 and 2), AtHRC1

Published Information

30 G1073 has been identified in the sequence of a BAC clone from chromosome 4 (BAC clone F23E12, gene F23E12.50, GenBank accession number AL022604), released by EU *Arabidopsis* Sequencing Project.

Closely Related Genes from Other Species

G1073 has similarity to *Medicago truncatula* cDNA clones (GenBank accession number AW574000 and AW560824) and *Glycine max* cDNA clones (AW349284 and AI736668) in the database.

Experimental Observations; Increased biomass and size, and other observations

The function of G1073 was analyzed using transgenic plants in which G1073 was expressed under the control of the cauliflower mosaic virus 35S promoter (these transgenic plants are referred to as "35S::G1073"). Transgenic plants overexpressing G1073 were substantially larger than wild-type controls, with at least a 60% increase in biomass (Figures 6A and 6B, 7A, and 7B; Table 6) . The increased mass of 35S::G1073 transgenic plants was attributed to enlargement of multiple organ types including stems, roots and floral organs; other than the size differences, these organs were not affected in their overall morphology. 35S::G1073 plants exhibited an increase of the width (but not length) of mature leaf organs, produced 2-3 more rosette leaves, and had enlarged cauline leaves in comparison to corresponding wild-type leaves. Overexpression of G1073 resulted in an increase in both leaf mass and leaf area per plant, and leaf morphology (G1073 overexpressors tended to produce more serrated leaves). We also found that root mass was increased in the transgenic plants, and that floral organs were also enlarged (Figures 7B). An increase of approximately 40% in stem diameter was observed in the transgenic plants. Images from the stem cross-sections of 35S::G1073 plants revealed that cortical cells are large and that vascular bundles contained more cells in the phloem and xylem relative to wild type (Figures 6A and 6B). Petal size in the 35S::G1073 lines was increased by 40-50% compared to wild type controls. Petal epidermal cells in those same lines were approximately 25-30% larger than those of the control plants. Furthermore, 15-20% more epidermal cells per petal were produced compared to wild type. Thus, in petals and stems, the increase in size was associated with an increase in cell size as well as in cell number.

Seed yield was also increased compared to control plants. 5S::G1073 lines showed an increase of at least 70% in seed yield (Table 6). This increased seed production was associated with an increased number of siliques per plant (Figure 10), rather than seeds per silique.

Table 6. Comparison of biomass and seed yield production in *Arabidopsis* wild-type and two 35S::G1073 overexpressing lines

Line	Fresh Weight (g)	Dry Weight (g)	Seed (g)
Wild-type	3.43 ± 0.70	0.73 ± 0.20	0.17 ± 0.07
35S::G1073-3	5.74 ± 1.74	1.17 ± 0.30	0.31 ± 0.08
35S::G1073-4	6.54 ± 2.19	1.38 ± 0.44	0.35 ± 0.12

All 35S::G1073 lines tested (10/10) exhibited significantly improved salt tolerance. Most of these lines also showed a sugar sensing phenotype, exhibiting improved germination on high sucrose media. One line showed increased heat germination tolerance. Flowering of G1073 overexpressing plants was delayed. Leaves of G1073 overexpressing plants were generally more serrated than those of wild-type plants. Improved drought tolerance was observed in 35S::G1073 transgenic lines.

A number of the CUT1::G1073 lines tested exhibited significantly improved salt tolerance and sugar sensing on high sucrose. One line showed improved germination on high mannitol.

Half of the ARSK::G1073 lines tested (5/10) showed improved germination on high salt, and two lines showed improved germination in cold relative to controls.

Utilities of G1073

Large size and late flowering produced as a result of G1073 or equivalent overexpression would be extremely useful in crops where the vegetative portion of the plant is the marketable portion (often vegetative growth stops when plants make the transition to flowering). In this case, it would be advantageous to prevent or delay flowering with the use of this gene or its equivalents in order to increase yield (biomass). Prevention of flowering by this gene or its equivalents would be useful in these same crops in order to prevent the spread of transgenic pollen and/or to prevent seed set. This gene or its equivalents could also be used to manipulate leaf shape, abiotic stress tolerance, including drought and salt tolerance, and seed yield.

G1067 (SEQ ID NOs: 3 and 4), AtHRC2

Published Information

A partial sequence of G1067 was identified from public EST clones (GenBank accession numbers W43561 and T43108). Weigel's group (The Salk Institute for Biological Studies) has recently identified an activation tagged mutant in which G1067 was overexpressed. The activation tagged mutant plants exhibited a late flowering phenotype in long days. Mutant leaves appeared wavy instead of flat, darker

green, larger, and rounder than those of wild type. Moreover, both leaf petioles and stem internodes were shorter than those of wild type (Weigel et al. (2000) *Plant Physiol.* 122:1003-1103).

Closely Related Genes from Other Species

5 G1067 is homologous to a *Medicago truncatula* cDNA clone (acc#AW574000).

Experimental observations

G1067 is a proprietary sequence discovered by us, and was initially identified from public EST clones (GenBank accession numbers W43561 and T43108). Full-length cDNA clones were later obtained
10 from our embryo specific cDNA library. The function of G1067 was analyzed using transgenic plants in which G1067 was expressed under the control of the 35S promoter.

A number of lines of transgenic plants overexpressing G1067 were found to be large and had broad leaves.

A number of different primary transformant lines of G1067 were also small with very twisted and
15 upcurled rosette leaves. In general these plants were poorly fertile, but sufficient seed was obtained from three plants for further analysis. Plants from these T2 lines were somewhat small with moderately curled leaves which had an undulating surface rather than the usual convex surface seen in wild-type leaves. One line with severely curled leaves also showed a lack of petiole extension reminiscent of the more severe phenotypes observed in the T1 generation. Biochemical analyses revealed that this line had low seed
20 protein.

G1067 appeared to be highly expressed in root and embryo. Its expression levels were also detected in siliques and germinating seeds. Expression of G1067 apparently is induced by auxin treatments.

ARSK1::G1067 overexpressing plants also showed increased tolerance in plate-based salt and
25 drought stress assays.

Utilities of G1067

Large size and late flowering produced as a result of G1067 or equivalent overexpression would be very useful for increasing vegetative portion of the plant This gene or its equivalents could also be used to
30 manipulate leaf shape or other aspects of plant architecture, and increase salt and drought tolerance.

G2153 (SEQ ID NOs: 5 and 6), AtHRC3

Published Information

The sequence of G2153 was obtained from *Arabidopsis* genomic sequencing project, GenBank accession number AC011437, based on its sequence similarity within the conserved domain to other AT-hook related proteins in *Arabidopsis*. G2153 corresponds to gene F7O18.4 (AAF04888). To date, there is no published information regarding the functions of this gene.

Closely Related Genes from Other Species

G2153 protein shows extensive sequence similarity with *Oryza sativa* chromosome 2 and 8 clones (AP004020 and AP003891), a *Lotus japonicus* cDNA (AW720668) and a *Medicago truncatula* cDNA clone (AW574000).

Experimental observations

The complete sequence of G2153 was determined by us. G2153 is strongly expressed in roots, embryos, siliques, and germinating seed, but at low or undetectable levels in shoots, flowers, and rosette leaves. It is not significantly induced or repressed by any condition tested.

The function of this gene was analyzed using transgenic plants in which G2153 was expressed under the control of the 35S promoter. A number of G2153 overexpressing lines were larger, and had broader, flatter leaves than those of wild-type plants. Some of these lines showed much larger rosettes than wild-type plants.

Overexpression of G2153 in *Arabidopsis* also resulted in seedlings with an altered response to osmotic stress. In a germination assay on media containing high sucrose, G2153 overexpressors had more expanded cotyledons and longer roots than the wild-type controls. This phenotype was confirmed in repeat experiments on individual lines, and all three lines showed osmotic tolerance. Increased tolerance to high sucrose could also be indicative of effects on sugar sensing. Overexpression of G2153 produced no consistent effects on *Arabidopsis* morphology, and no altered phenotypes were noted in any of the biochemical assays.

G2153 was also overexpressed in tomato plants that were then used in field trials. At one stage in the trial, the plants were deprived of water for several days. Upon subsequent watering, a number of the transgenic plants were found to be larger and healthier than wild-type tomato plants, and at least one line produced more fruit than wild-type plants.

Utilities of G2153

G2153 could be used to increase a plant's biomass.

G2153 may be useful for altering a plant's response to sugars, and may also be used to alter a plant's response to water deficit conditions. Therefore, G2153 could be used to engineer plants with enhanced tolerance to drought, salt stress, and freezing.

G2156 (SEQ ID NOs: 7 and 8), AtHRC4

Published Information

The sequence of G2156 was obtained from *Arabidopsis* genomic sequencing project, GenBank accession number AC015450, based on its sequence similarity within the conserved domain to other AT-hook related proteins in *Arabidopsis*. G2156 corresponds to gene F14G6.10 (AAG51949). To date, there is no published information regarding the functions of this gene.

Closely Related Genes from Other Species

G2156 protein shows extensive sequence similarity with *Medicago truncatula* cDNA clones (AW574000 and AW774484) and a *Lycopersicon esculentum* cDNA clone (BG134451).

Experimental Observations

The complete sequence of G2156 was determined by us. G2156 was found to be expressed at moderate levels in embryos and siliques, and at significantly lower levels in roots, flowers, and germinating seed. It shows possible induction by auxin.

The function of this gene was analyzed using transgenic plants in which G2156 was expressed under the control of the 35S promoter. A majority (8 of 10) of the 35S::G2156 transformants tested showed tolerance to high salt concentrations in plate-based assays. One line also showed a strong sugar-sensing phenotype. Another line showed tolerance to germination in heat.

The function of this gene was also analyzed using transgenic plants in which the gene was expressed under the control of the ARSK1 promoter. ARSK1::G2156 overexpressing plants were shown to be more drought tolerant than wild-type control plants in soil-based assays.

A number of *Arabidopsis* lines overexpressing G2156 under the control of the 35S promoter were found to be larger, with broader leaves and larger rosettes than wild-type control plants.

Utilities of G2156

G2156 could be used to increase a plant's biomass.

G2156 could be used to improve a plant's germination in hot conditions, and also improve cold tolerance.

5 G2156 could be also used to alter a plant's response to water deficit conditions and, therefore, could be used to engineer plants with enhanced tolerance to drought, salt stress, and freezing.

G2153 may also be useful for altering a plant's response to sugars.

Rice sequences G3399 and G3407 (SEQ ID NOs: 9-12), OsHRC2 and OsHRC7

10 Published Information

The sequences of G3399 and G3407 were discovered based on their similarity to G1073 as determined by BLAST analysis of a proprietary database , To date, there is no published information regarding the functions of either gene or polypeptide.

15 Experimental Observations

A number of *Arabidopsis* lines overexpressing G3399 and G3407 under the control of the 35S promoter were found be larger, with broader leaves and larger rosettes than wild-type control plants.

Utilities of G3399 and G3407

20 G3399 and G3407 could be used to increase a plant's biomass.

G3399 and G3407 may be also used to alter a plant's response to water deficit conditions and, therefore, could be used to engineer plants with enhanced tolerance to drought, salt stress, and freezing.

Soybean sequences G3456,G3459 and G3460 (SEQ ID NOs: 13-18), GmHRC2, GmHRC7 and

25 **GmHRC8**

Published Information

The sequences of G3456,G3459 and G3460 were discovered based on their similarity to G1073 as determined by BLAST analysis of a proprietary database , To date, there is no published information regarding the functions of either gene or polypeptide.

30

Experimental Observations

A significant number of *Arabidopsis* lines overexpressing G3456,G3459 and G3460 under the control of the 35S promoter were found be larger, with broader leaves and larger rosettes than wild-type

control plants.

Utilities of G3456, G3459 and G3460

G3456, G3459 and G3460 can be used to increase a plant's biomass.

- 5 G3456, G3459 and G3460 may be also used to alter a plant's response to water deficit conditions and, therefore, could be used to engineer plants with enhanced tolerance to drought, salt stress, and freezing.

Example IX: Identification of Homologous Sequences

- 10 This example describes identification of genes that are orthologous to *Arabidopsis thaliana* transcription factors from a computer homology search.

- Homologous sequences, including those of paralogs and orthologs from *Arabidopsis* and other plant species, were identified using database sequence search tools, such as the Basic Local Alignment Search Tool (BLAST) (Altschul et al. (1990) *J. Mol. Biol.* 215: 403-410; and Altschul et al. (1997) *Nucleic Acid Res.* 25: 3389-3402). The tblastx sequence analysis programs were employed using the BLOSUM-62 scoring matrix (Henikoff and Henikoff (1992) *Proc. Natl. Acad. Sci. USA* 89: 10915-10919). The entire NCBI GenBank database was filtered for sequences from all plants except *Arabidopsis thaliana* by selecting all entries in the NCBI GenBank database associated with NCBI taxonomic ID 33090 (Viridiplantae; all plants) and excluding entries associated with taxonomic ID 3701 (*Arabidopsis thaliana*).
- 20

- These sequences are compared to sequences representing transcription factor genes presented in the Sequence Listing, using the Washington University TBLASTX algorithm (version 2.0a19MP) at the default settings using gapped alignments with the filter "off". For each transcription factor gene in the Sequence Listing, individual comparisons were ordered by probability score (P-value), where the score reflects the probability that a particular alignment occurred by chance. For example, a score of $3.6e-59$ is 3.6×10^{-59} . In addition to P-values, comparisons were also scored by percentage identity. Percentage identity reflects the degree to which two segments of DNA or protein are identical over a particular length. Examples of sequences so identified are presented in, for example, the Sequence Listing, and Table 5. Paralogous or orthologous sequences were readily identified and available in GenBank by Accession number (Table 5; Sequence Identifier or Accession Number). The percent sequence identity among these sequences can be as low as 49%, or even lower sequence identity.
- 25
- 30

Candidate paralogous sequences were identified among *Arabidopsis* transcription factors through alignment, identity, and phylogenetic relationships. G1067, G2153 and G2156 (SEQ ID NO: 4, 6, and 8, respectively), paralogs of G1073, may be found in the Sequence Listing.

Candidate orthologous sequences were identified from proprietary unigene sets of plant gene sequences in *Zea mays*, *Glycine max* and *Oryza sativa* based on significant homology to *Arabidopsis* transcription factors. These candidates were reciprocally compared to the set of *Arabidopsis* transcription factors. If the candidate showed maximal similarity in the protein domain to the eliciting transcription factor or to a paralog of the eliciting transcription factor, then it was considered to be an ortholog. Identified non-*Arabidopsis* sequences that were shown in this manner to be orthologous to the *Arabidopsis* sequences are provided in, for example, Table 5.

Example X: Screen of Plant cDNA library for Sequence Encoding a Transcription Factor DNA Binding Domain That Binds To a Transcription Factor Binding Promoter Element and Demonstration of Protein Transcription Regulation Activity.

The "one-hybrid" strategy (Li and Herskowitz (1993) *Science* 262: 1870-1874) is used to screen for plant cDNA clones encoding a polypeptide comprising a transcription factor DNA binding domain, a conserved domain. In brief, yeast strains are constructed that contain a lacZ reporter gene with either wild-type or mutant transcription factor binding promoter element sequences in place of the normal UAS (upstream activator sequence) of the GALL promoter. Yeast reporter strains are constructed that carry transcription factor binding promoter element sequences as UAS elements are operably linked upstream (5') of a lacZ reporter gene with a minimal GAL1 promoter. The strains are transformed with a plant expression library that contains random cDNA inserts fused to the GAL4 activation domain (GAL4-ACT) and screened for blue colony formation on X-gal-treated filters (X-gal: 5-bromo-4-chloro-3-indolyl- β -D-galactoside; Invitrogen Corporation, Carlsbad CA). Alternatively, the strains are transformed with a cDNA polynucleotide encoding a known transcription factor DNA binding domain polypeptide sequence.

Yeast strains carrying these reporter constructs produce low levels of beta-galactosidase and form white colonies on filters containing X-gal. The reporter strains carrying wild-type transcription factor binding promoter element sequences are transformed with a polynucleotide that encodes a polypeptide comprising a plant transcription factor DNA binding domain operably linked to the acidic activator domain of the yeast GAL4 transcription factor, "GAL4-ACT". The clones that contain a polynucleotide encoding a transcription factor DNA binding domain operably linked to GLA4-ACT can bind upstream of the lacZ reporter genes carrying the wild-type transcription factor binding promoter element sequence, activate transcription of the lacZ gene and result in yeast forming blue colonies on X-gal-treated filters.

Upon screening about 2×10^6 yeast transformants, positive cDNA clones are isolated; i.e., clones that cause yeast strains carrying lacZ reporters operably linked to wild-type transcription factor binding promoter elements to form blue colonies on X-gal-treated filters. The cDNA clones do not cause a yeast strain carrying a mutant type transcription factor binding promoter elements fused to LacZ to turn blue.

- 5 Thus, a polynucleotide encoding transcription factor DNA binding domain, a conserved domain, is shown to activate transcription of a gene.

Example XI: Gel Shift Assays.

- The presence of a transcription factor comprising a DNA binding domain which binds to a DNA transcription factor binding element is evaluated using the following gel shift assay. The transcription factor is recombinantly expressed and isolated from *E. coli* or isolated from plant material. Total soluble protein, including transcription factor, (40 ng) is incubated at room temperature in 10 μ l of 1 x binding buffer (15 mM HEPES (pH 7.9), 1 mM EDTA, 30 mM KCl, 5% glycerol, 5% bovine serum albumin, 1 mM DTT) plus 50 ng poly(dI-dC):poly(dI-dC) (Pharmacia, Piscataway NJ) with or without 100 ng competitor DNA. After 10 minutes incubation, probe DNA comprising a DNA transcription factor binding element (1 ng) that has been 32 P-labeled by end-filling (Sambrook et al. (1989) *supra*) is added and the mixture incubated for an additional 10 minutes. Samples are loaded onto polyacrylamide gels (4% w/v) and fractionated by electrophoresis at 150V for 2h (Sambrook et al. *supra*). The degree of transcription factor-probe DNA binding is visualized using autoradiography. Probes and competitor DNAs are prepared from oligonucleotide inserts ligated into the BamHI site of pUC118 (Vieira et al. (1987) *Methods Enzymol.* 153: 3-11). Orientation and concatenation number of the inserts are determined by dideoxy DNA sequence analysis (Sambrook et al. *supra*). Inserts are recovered after restriction digestion with EcoRI and HindIII and fractionation on polyacrylamide gels (12% w/v) (Sambrook et al. *supra*).

Example XII. Introduction of Polynucleotides into Dicotyledonous Plants

- Transcription factor sequences listed in the Sequence Listing recombined into pMEN20 or pMEN65 expression vectors are transformed into a plant for the purpose of modifying plant traits. The cloning vector may be introduced into a variety of cereal plants by means well known in the art such as, for example, direct DNA transfer or *Agrobacterium tumefaciens*-mediated transformation. It is now routine to produce transgenic plants using most dicot plants (see Weissbach and Weissbach, (1989) *supra*; Gelvin et al. (1990) *supra*; Herrera-Estrella et al. (1983) *supra*; Bevan (1984) *supra*; and Klee (1985) *supra*). Methods for analysis of traits are routine in the art and examples are disclosed above.

Example XIII: Transformation of Cereal Plants with an Expression Vector

Cereal plants such as, but not limited to, corn, wheat, rice, sorghum, or barley, may also be transformed with the present polynucleotide sequences in pMEN20 or pMEN65 expression vectors for the purpose of modifying plant traits. For example, pMEN020 may be modified to replace the NptII coding region with the BAR gene of *Streptomyces hygrosopicus* that confers resistance to phosphinothricin. The KpnI and BglII sites of the Bar gene are removed by site-directed mutagenesis with silent codon changes.

The cloning vector may be introduced into a variety of cereal plants by means well known in the art such as, for example, direct DNA transfer or *Agrobacterium tumefaciens*-mediated transformation. It is now routine to produce transgenic plants of most cereal crops (Vasil (1994) *Plant Mol. Biol.* 25: 925-937) such as corn, wheat, rice, sorghum (Cassas et al. (1993) *Proc. Natl. Acad. Sci.* 90: 11212-11216, and barley (Wan and Lemeaux (1994) *Plant Physiol.* 104:37-48. DNA transfer methods such as the microprojectile can be used for corn (Fromm et al. (1990) *Bio/Technol.* 8: 833-839); Gordon-Kamm et al. (1990) *Plant Cell* 2: 603-618; Ishida (1990) *Nature Biotechnol.* 14:745-750), wheat (Vasil et al. (1992) *Bio/Technol.* 10:667-674; Vasil et al. (1993) *Bio/Technol.* 11:1553-1558; Weeks et al. (1993) *Plant Physiol.* 102:1077-1084), rice (Christou (1991) *Bio/Technol.* 9:957-962; Hiei et al. (1994) *Plant J.* 6:271-282; Aldemita and Hodges (1996) *Planta* 199:612-617; and Hiei et al. (1997) *Plant Mol. Biol.* 35:205-218). For most cereal plants, embryogenic cells derived from immature scutellum tissues are the preferred cellular targets for transformation (Hiei et al. (1997) *Plant Mol. Biol.* 35:205-218; Vasil (1994) *Plant Mol. Biol.* 25: 925-937).

Vectors according to the present invention may be transformed into corn embryogenic cells derived from immature scutellar tissue by using microprojectile bombardment, with the A188XB73 genotype as the preferred genotype (Fromm et al. (1990) *Bio/Technol.* 8: 833-839; Gordon-Kamm et al. (1990) *Plant Cell* 2: 603-618). After microprojectile bombardment the tissues are selected on phosphinothricin to identify the transgenic embryogenic cells (Gordon-Kamm et al. (1990) *Plant Cell* 2: 603-618). Transgenic plants are regenerated by standard corn regeneration techniques (Fromm et al. (1990) *Bio/Technol.* 8: 833-839; Gordon-Kamm et al. (1990) *Plant Cell* 2: 603-618).

The plasmids prepared as described above can also be used to produce transgenic wheat and rice plants (Christou (1991) *Bio/Technol.* 9:957-962; Hiei et al. (1994) *Plant J.* 6:271-282; Aldemita and Hodges (1996) *Planta* 199:612-617; and Hiei et al. (1997) *Plant Mol. Biol.* 35:205-218) that coordinately express genes of interest by following standard transformation protocols known to those skilled in the art for rice and wheat (Vasil et al. (1992) *Bio/Technol.* 10:667-674; Vasil et al. (1993) *Bio/Technol.* 11:1553-1558; and Weeks et al. (1993) *Plant Physiol.* 102:1077-1084), where the bar gene is used as the selectable marker.

Example XIV: Transformation of Tomato and Soy Plants

Numerous protocols for the transformation of tomato and soy plants have been previously described, and are well known in the art. Gruber et al. ((1993) in Methods in Plant Molecular Biology and Biotechnology, p. 89-119, Glick and Thompson, eds., CRC Press, Inc., Boca Raton) describe several expression vectors and culture methods that may be used for cell or tissue transformation and subsequent regeneration. For soybean transformation, methods are described by Miki et al. (1993) in Methods in Plant Molecular Biology and Biotechnology, p. 67-88, Glick and Thompson, eds., CRC Press, Inc., Boca Raton; and U.S. Pat. No. 5,563,055, (Townsend and Thomas), issued Oct.8, 1996.

There are a substantial number of alternatives to *Agrobacterium*-mediated transformation protocols, other methods for the purpose of transferring exogenous genes into soybeans or tomatoes. One such method is microprojectile-mediated transformation, in which DNA on the surface of microprojectile particles is driven into plant tissues with a biolistic device (see, for example, Sanford et al., (1987) *Part. Sci. Technol.* 5:27-37; Christou et al. (1992) *Plant. J.* 2: 275-281; Sanford (1993) *Methods Enzymol.* 217: 483-509; Klein et al. (1987) *Nature* 327: 70-73; U.S. Pat. No.5,015,580 (Christou et al), issued May 14, 1991; and U.S. Pat. No. 5,322,783 (Tomes et al.), issued Jun. 21, 1994.

Alternatively, sonication methods (see, for example, Zhang et al. (1991) *Bio/Technology* 9: 996-997); direct uptake of DNA into protoplasts using CaCl₂ precipitation, polyvinyl alcohol or poly-L-ornithine (see, for example, Hain et al. (1985) *Mol. Gen. Genet.* 199: 161-168; Draper et al., *Plant Cell Physiol.* 23: 451-458 (1982)); liposome or spheroplast fusion (see, for example, Deshayes et al. (1985) *EMBO J.*, 4: 2731-2737; Christou et al. (1987) *Proc. Natl. Acad. Sci. U.S.A.* 84: 3962-3966); and electroporation of protoplasts and whole cells and tissues (see, for example, Donn et al.(1990) in Abstracts of VIIth International Congress on Plant Cell and Tissue Culture IAPTC, A2-38: 53; D'Halluin et al. (1992) *Plant Cell* 4: 1495-1505;and Spencer et al. (1994) *Plant Mol. Biol.* 24: 51-61) have been used to introduce foreign DNA and expression vectors into plants.

After plants or plant cells are transformed (and the latter regenerated into plants) the transgenic plant thus generated may be crossed with itself or a plant from the same line, a non-transformed or wild-type plant, or another transformed plant from a different transgenic line of plants. Crossing provides the advantages of being able to produce new and perhaps stable transgenic varieties. Genes and the traits they confer that have been introduced into a tomato or soybean line may be moved into distinct line of plants using traditional backcrossing techniques well known in the art. Transformation of tomato plants may be conducted using the protocols of Koornneef et al (1986) In Tomato Biotechnology: Alan R. Liss, Inc., 169-178, and in U.S. Patent 6,613,962, the latter method described in brief here. Eight day old cotyledon explants are precultured for 24 hours in Petri dishes containing a feeder layer of *Petunia hybrida*

suspension cells plated on MS medium with 2% (w/v) sucrose and 0.8% agar supplemented with 10 μ M α -naphthalene acetic acid and 4.4 μ M 6-benzylaminopurine. The explants are then infected with a diluted overnight culture of *Agrobacterium tumefaciens* containing an expression vector comprising a polynucleotide of the invention for 5-10 minutes, blotted dry on sterile filter paper and cocultured for 48 hours on the original feeder layer plates. Culture conditions are as described above. Overnight cultures of *Agrobacterium tumefaciens* are diluted in liquid MS medium with 2% (w/v) sucrose, pH 5.7) to an OD₆₀₀ of 0.8.

Following the cocultivation, the cotyledon explants are transferred to Petri dishes with selective medium consisting of MS medium supplemented with 4.56 μ M zeatin, 67.3 μ M vancomycin, 418.9 μ M cefotaxime and 171.6 μ M kanamycin sulfate, and cultured under the culture conditions described above. The explants are subcultured every three weeks onto fresh medium. Emerging shoots are dissected from the underlying callus and transferred to glass jars with selective medium without zeatin to form roots. The formation of roots in a medium containing kanamycin sulphate is regarded as a positive indication of a successful transformation.

Transformation of soybean plants may be conducted using the methods found in, for example, U.S. Patent 5,563,055 (Townsend et al., issued October 8, 1996), described in brief here. In this method soybean seed is surface sterilized by exposure to chlorine gas evolved in a glass bell jar. Seeds are germinated by plating on 1/10 strength agar solidified medium without plant growth regulators and culturing at 28° C. with a 16 hour day length. After three or four days, seed may be prepared for cocultivation. The seedcoat is removed and the elongating radicle removed 3-4 mm below the cotyledons.

Overnight cultures of *Agrobacterium tumefaciens* harboring the expression vector comprising a polynucleotide of the invention are grown to log phase, pooled, and concentrated by centrifugation. Inoculations are conducted in batches such that each plate of seed was treated with a newly resuspended pellet of *Agrobacterium*. The pellets are resuspended in 20 ml inoculation medium. The inoculum is poured into a Petri dish containing prepared seed and the cotyledonary nodes are macerated with a surgical blade. After 30 minutes the explants are transferred to plates of the same medium which has been solidified. Explants are embedded with the adaxial side up and level with the surface of the medium and cultured at 22° C. for three days under white fluorescent light. These plants may then be regenerated according to methods well established in the art, such as by moving the explants after three days to a liquid counter-selection medium (see U.S. Patent 5,563,055).

The explants may then be picked, embedded and cultured in solidified selection medium. After one month on selective media transformed tissue becomes visible as green sectors of regenerating tissue against a background of bleached, less healthy tissue. Explants with green sectors are transferred to an

elongation medium. Culture is continued on this medium with transfers to fresh plates every two weeks. When shoots are 0.5 cm in length they may be excised at the base and placed in a rooting medium.

Example XV: Genes that Confer Significant Improvements to non-*Arabidopsis* species

5 The function of specific orthologs of G1073 have been analyzed and may be further characterized through their ectopic overexpression in plants, using the CaMV 35S, ARSK1, or other appropriate promoter, identified above. Genes that have been examined and have been shown to modify plant traits (including increasing biomass and abiotic stress tolerance) encode members of the AT-hook transcription factors, such as those found in *Arabidopsis thaliana* (SEQ ID NO: 2, 4, 6 and 8) *Oryza sativa* (SEQ ID
10 NO: 10 and 12), and *Glycine max* (SEQ ID NO: 14, 16 and 18). In addition to these sequences, it is expected that related polynucleotide sequences encoding polypeptides found the Sequence Listing can also induce altered traits, including increased biomass and abiotic stress tolerance, when transformed into a variety of plants. The polynucleotide and polypeptide sequences derived from monocots (e.g., the rice sequences) may be used to transform both monocot and dicot plants, and those derived from dicots (e.g.,
15 the *Arabidopsis* and soy genes) may be used to transform either group, although some of these sequences will function best if the gene is transformed into a plant from the same group as that from which the sequence is derived.

Seeds of these transgenic plants are subjected to germination assays to measure sucrose sensing. Sterile monocot seeds, including, but not limited to, corn, rice, wheat, rye and sorghum, as well as dicots
20 including, but not limited to soybean and alfalfa, are sown on 80% MS medium plus vitamins with 9.4% sucrose; control media lack sucrose. All assay plates are then incubated at 22° C under 24-hour light, 120-130 $\mu\text{Ein}/\text{m}^2/\text{s}$, in a growth chamber. Evaluation of germination and seedling vigor is then conducted three days after planting. Overexpressors of these genes may be found to be more tolerant to high sucrose by having better germination, longer radicles, and more cotyledon expansion. These results have previously
25 indicated that overexpressors of G1073, G1067, G2153 and/or G2156 are involved in sucrose-specific sugar sensing; it is expected that structurally similar orthologs of these sequences, including those found in the Sequence Listing, are also be involved in sugar sensing, an indication of altered osmotic stress tolerance.

Plants overexpressing these orthologs may also be subjected to soil-based drought assays to
30 identify those lines that are more tolerant to water deprivation than wild-type control plants. Generally, 35S:: or ARSK1::G1073, G1067, G2153 and/or G2156 ortholog overexpressing plants will appear significantly larger and greener, with less wilting or desiccation, than wild-type controls plants, particularly after a period of water deprivation is followed by rewatering and a subsequent incubation period.

Monocotyledonous plants such as rice, corn, wheat, rye, sorghum, barley and others may be transformed with a plasmid containing G1073, G1067, G2153, G2156, G3399, G3407, G3456, G3459 and G3460 equivalents, including monocot-derived sequences such as those presented in Table 5, or AT-hook transcription factor genes, cloned into a vector such as pGA643 and containing a kanamycin-resistance marker, and are expressed constitutively under the CaMV 35S promoter or COR15 promoter.

The cloning vector may be introduced into monocots by, for example, means described in detail in Example XIII, including direct DNA transfer or *Agrobacterium tumefaciens*-mediated transformation. The latter approach may be accomplished by a variety of means, including, for example, that of U.S. Patent No. 5,591,616, in which monocotyledon callus is transformed by contacting dedifferentiating tissue with the *Agrobacterium* containing the cloning vector.

The sample tissues are immersed in a suspension of 3×10^9 cells of *Agrobacterium* containing the cloning vector for 3-10 minutes. The callus material is cultured on solid medium at 25° C in the dark for several days. The calli grown on this medium are transferred to Regeneration medium. Transfers are continued every 2-3 weeks (2 or 3 times) until shoots develop. Shoots are then transferred to Shoot-Elongation medium every 2-3 weeks. Healthy looking shoots are transferred to rooting medium and after roots have developed, the plants are placed into moist potting soil.

The transformed plants are then analyzed for the presence of the NPTII gene/ kanamycin resistance by ELISA, using the ELISA NPTII kit from 5Prime-3Prime Inc. (Boulder, CO).

Northern blot analysis, RT-PCR or microarray analysis of the regenerated, transformed plants may be used to show expression of G1073, G1067, G2153, G2156, G3399, G3407, G3456, G3459 and G3460 equivalent genes that are capable of inducing abiotic stress tolerance.

To verify the ability to confer abiotic stress tolerance, mature plants expressing a monocot-derived equivalent gene, or alternatively, seedling progeny of these plants, may be challenged using stresses described in Example XV. By comparing wild type plants and the transgenic plants, the latter are shown to be more tolerant to abiotic stress, and/or have increased biomass, as compared to wild type control plant similarly treated.

These experiments demonstrate that equivalents of G1073, G1067, G2153, G2156, G3399, G3407, G3456, G3459 and G3460 can be identified and shown to increase biomass and improve abiotic stress tolerance, including osmotic stresses such as drought or salt stress.

Example XVI: Identification of Orthologous and Paralogous Sequences by PCR

Orthologs to *Arabidopsis* genes may be identified by several methods, including hybridization, amplification, or bioinformatically. This example describes how one may identify equivalents to the

Arabidopsis AP2 family transcription factor CBF1 (polynucleotide SEQ ID NO: 69, encoded polypeptide SEQ ID NO: 70), which confers tolerance to abiotic stresses (Thomashow et al. (2002) U.S. Patent No. 6,417,428), and an example to confirm the function of homologous sequences. In this example, orthologs to CBF1 were found in canola (*Brassica napus*) using polymerase chain reaction (PCR).

Degenerate primers were designed for regions of AP2 binding domain and outside of the AP2 (carboxyl terminal domain):

Mol 368 (reverse) 5'- CAY CCN ATH TAY MGN GGN GT -3' (SEQ ID NO: 77)

Mol 378 (forward) 5'- GGN ARN ARC ATN CCY TCN GCC -3' (SEQ ID NO: 78)

(Y: C/T, N: A/C/G/T, H: A/C/T, M: A/C, R: A/G)

Primer Mol 368 is in the AP2 binding domain of CBF1 (amino acid sequence: His-Pro-Ile-Tyr-Arg-Gly-Val) while primer Mol 378 is outside the AP2 domain (carboxyl terminal domain) (amino acid sequence: Met-Ala-Glu-Gly-Met-Leu-Leu-Pro).

The genomic DNA isolated from *B. napus* was PCR-amplified by using these primers following these conditions: an initial denaturation step of 2 min at 93° C; 35 cycles of 93° C for 1 min, 55° C for 1 min, and 72° C for 1 min ; and a final incubation of 7 min at 72° C at the end of cycling.

The PCR products were separated by electrophoresis on a 1.2% agarose gel and transferred to nylon membrane and hybridized with the AT CBF1 probe prepared from *Arabidopsis* genomic DNA by PCR amplification. The hybridized products were visualized by colorimetric detection system (Boehringer Mannheim) and the corresponding bands from a similar agarose gel were isolated using the Qiagen Extraction Kit (Qiagen, Valencia CA). The DNA fragments were ligated into the TA clone vector from TOPO TA Cloning Kit (Invitrogen Corporation, Carlsbad CA) and transformed into *E. coli* strain TOP10 (Invitrogen).

Seven colonies were picked and the inserts were sequenced on an ABI 377 machine from both strands of sense and antisense after plasmid DNA isolation. The DNA sequence was edited by sequencer and aligned with the AtCBF1 by GCG software and NCBI blast searching.

The nucleic acid sequence and amino acid sequence of one canola ortholog found in this manner (bnCBF1; polynucleotide SEQ ID NO: 75 and polypeptide SEQ ID NO: 76) identified by this process is shown in the Sequence Listing.

The aligned amino acid sequences show that the bnCBF1 gene has 88% identity with the *Arabidopsis* sequence in the AP2 domain region and 85% identity with the *Arabidopsis* sequence outside the AP2 domain when aligned for two insertion sequences that are outside the AP2 domain.

Similarly, paralogous sequences to *Arabidopsis* genes, such as *CBF1*, may also be identified.

Two paralogs of CBF1 from *Arabidopsis thaliana*: *CBF2* and *CBF3*. *CBF2* and *CBF3* have been cloned and sequenced as described below. The sequences of the DNA SEQ ID NO: 71 and 73 and encoded proteins SEQ ID NO: 72 and 74 are set forth in the Sequence Listing.

A lambda cDNA library prepared from RNA isolated from *Arabidopsis thaliana* ecotype Columbia (Lin and Thomashow (1992) *Plant Physiol.* 99: 519-525) was screened for recombinant clones that carried inserts related to the *CBF1* gene (Stockinger et al. (1997) *Proc. Natl. Acad. Sci.* 94:1035-1040). CBF1 was ³²P-radiolabeled by random priming (Sambrook et al. *supra*) and used to screen the library by the plaque-lift technique using standard stringent hybridization and wash conditions (Hajela et al. (1990) *Plant Physiol.* 93:1246-1252; Sambrook et al. *supra*) 6 X SSPE buffer, 60° C for hybridization and 0.1 X SSPE buffer and 60° C for washes). Twelve positively hybridizing clones were obtained and the DNA sequences of the cDNA inserts were determined. The results indicated that the clones fell into three classes. One class carried inserts corresponding to *CBF1*. The two other classes carried sequences corresponding to two different homologs of *CBF1*, designated *CBF2* and *CBF3*. The nucleic acid sequences and predicted protein coding sequences for *Arabidopsis CBF1*, *CBF2* and *CBF3* are listed in the Sequence Listing (SEQ ID NOs: 69, 71, 73 and SEQ ID NOs: 70, 72, and 74, respectively). The nucleic acid sequences and predicted protein coding sequence for *Brassica napus* CBF ortholog is listed in the Sequence Listing (SEQ ID NOs: 75 and 76, respectively).

A comparison of the nucleic acid sequences of *Arabidopsis CBF1*, *CBF2* and *CBF3* indicate that they are 83 to 85% identical as shown in Table 7.

TABLE 7

	Percent identity ^a	
	DNA ^b	Polypeptide
cbf1/cbf2	85	86
cbf1/cbf3	83	84
cbf2/cbf3	84	85

^a Percent identity was determined using the *Clustal* algorithm from the Megalign program (DNASTAR, Inc.).

^b Comparisons of the nucleic acid sequences of the open reading frames are shown.

Similarly, the amino acid sequences of the three CBF polypeptides range from 84 to 86% identity. An alignment of the three amino acid sequences reveals that most of the differences in amino acid sequence occur in the acidic C-terminal half of the polypeptide. This region of CBF1 serves as an activation domain in both yeast and *Arabidopsis* (not shown).

Residues 47 to 106 of CBF1 correspond to the AP2 domain of the protein, a DNA binding motif that to date, has only been found in plant proteins. A comparison of the AP2 domains of CBF1, CBF2 and CBF3 indicates that there are a few differences in amino acid sequence. These differences in amino acid sequence might have an effect on DNA binding specificity.

Example XVII: Transformation of Canola with a Plasmid Containing CBF1, CBF2, or CBF3

After identifying homologous genes to CBF1, canola was transformed with a plasmid containing the *Arabidopsis* CBF1, CBF2, or CBF3 genes cloned into the vector pGA643 (An (1987) *Methods Enzymol.* 253: 292). In these constructs the CBF genes were expressed constitutively under the CaMV 35S promoter. In addition, the CBF1 gene was cloned under the control of the *Arabidopsis* COR15 promoter in the same vector pGA643. Each construct was transformed into *Agrobacterium* strain GV3101. Transformed *Agrobacteria* were grown for 2 days in minimal AB medium containing appropriate antibiotics.

Spring canola (*B. napus* cv. Westar) was transformed using the protocol of Moloney et al. (1989) *Plant Cell Reports* 8: 238, with some modifications as described. Briefly, seeds were sterilized and plated on half strength MS medium, containing 1% sucrose. Plates were incubated at 24° C under 60-80 $\mu\text{E}/\text{m}^2\text{s}$ light using a 16 hour light/ 8 hour dark photoperiod. Cotyledons from 4-5 day old seedlings were collected, the petioles cut and dipped into the *Agrobacterium* solution. The dipped cotyledons were placed on co-cultivation medium at a density of 20 cotyledons/plate and incubated as described above for 3 days. Explants were transferred to the same media, but containing 300 mg/l timentin (SmithKline Beecham, PA) and thinned to 10 cotyledons/plate. After 7 days explants were transferred to Selection/Regeneration medium. Transfers were continued every 2-3 weeks (2 or 3 times) until shoots had developed. Shoots were transferred to Shoot-Elongation medium every 2-3 weeks. Healthy looking shoots were transferred to rooting medium. Once good roots had developed, the plants were placed into moist potting soil.

The transformed plants were then analyzed for the presence of the NPTII gene/ kanamycin resistance by ELISA, using the ELISA NPTII kit from 5Prime-3Prime Inc. (Boulder, CO). Approximately 70% of the screened plants were NPTII positive. Only those plants were further analyzed.

From Northern blot analysis of the plants that were transformed with the constitutively expressing constructs, showed expression of the CBF genes and all CBF genes were capable of inducing the *Brassica napus* cold-regulated gene BN115 (homolog of the *Arabidopsis* COR15 gene). Most of the transgenic plants appear to exhibit a normal growth phenotype. As expected, the transgenic plants are more freezing tolerant than the wild-type plants. Using the electrolyte leakage of leaves test, the control showed a 50% leakage at -2 to -3° C. Spring canola transformed with either CBF1 or CBF2 showed a 50% leakage at -6 to -7° C. Spring canola transformed with CBF3 shows a 50% leakage at about -10 to -15° C. Winter canola transformed with CBF3 may show a 50% leakage at about -16 to -20° C. Furthermore, if the spring or winter canola are cold acclimated the transformed plants may exhibit a further increase in freezing tolerance of at least -2° C.

To test salinity tolerance of the transformed plants, plants were watered with 150 mM NaCl. Plants overexpressing CBF1, CBF2, or CBF3 grew better compared with plants that had not been transformed with CBF1, CBF2, or CBF3.

These results demonstrate that equivalents of *Arabidopsis* transcription factors can be identified and shown to confer similar functions in non-*Arabidopsis* plant species.

Example XVIII: Cloning of transcription factor promoters

Promoters are isolated from transcription factor genes that have gene expression patterns useful for a range of applications, as determined by methods well known in the art (including transcript profile analysis with cDNA or oligonucleotide microarrays, Northern blot analysis, semi-quantitative or quantitative RT-PCR). Interesting gene expression profiles are revealed by determining transcript abundance for a selected transcription factor gene after exposure of plants to a range of different experimental conditions, and in a range of different tissue or organ types, or developmental stages. Experimental conditions to which plants are exposed for this purpose includes cold, heat, drought, osmotic challenge, varied hormone concentrations (ABA, GA, auxin, cytokinin, salicylic acid, brassinosteroid), pathogen and pest challenge. The tissue types and developmental stages include stem, root, flower, rosette leaves, cauline leaves, siliques, germinating seed, and meristematic tissue. The set of expression levels provides a pattern that is determined by the regulatory elements of the gene promoter.

Transcription factor promoters for the genes disclosed herein are obtained by cloning 1.5 kb to 2.0 kb of genomic sequence immediately upstream of the translation start codon for the coding sequence of the encoded transcription factor protein. This region includes the 5'-UTR of the transcription factor gene, which can comprise regulatory elements. The 1.5 kb to 2.0 kb region is cloned through PCR methods, using primers that include one in the 3' direction located at the translation start codon (including

appropriate adaptor sequence), and one in the 5' direction located from 1.5 kb to 2.0 kb upstream of the translation start codon (including appropriate adaptor sequence). The desired fragments are PCR-amplified from *Arabidopsis* Col-0 genomic DNA using high-fidelity Taq DNA polymerase to minimize the incorporation of point mutation(s). The cloning primers incorporate two rare restriction sites, such as Not1 and Sfi1, found at low frequency throughout the *Arabidopsis* genome. Additional restriction sites are used in the instances where a Not1 or Sfi1 restriction site is present within the promoter.

The 1.5-2.0 kb fragment upstream from the translation start codon, including the 5'-untranslated region of the transcription factor, is cloned in a binary transformation vector immediately upstream of a suitable reporter gene, or a transactivator gene that is capable of programming expression of a reporter gene in a second gene construct. Reporter genes used include green fluorescent protein (and related fluorescent protein color variants), beta-glucuronidase, and luciferase. Suitable transactivator genes include LexA-GAL4, along with a transactivatable reporter in a second binary plasmid (as disclosed in U.S. patent application 09/958,131, incorporated herein by reference). The binary plasmid(s) is transferred into *Agrobacterium* and the structure of the plasmid confirmed by PCR. These strains are introduced into *Arabidopsis* plants as described in other examples, and gene expression patterns determined according to standard methods known to one skilled in the art for monitoring GFP fluorescence, beta-glucuronidase activity, or luminescence.

All publications and patent applications mentioned in this specification are herein incorporated by reference to the same extent as if each individual publication or patent application was specifically and individually indicated to be incorporated by reference.

The present invention is not limited by the specific embodiments described herein. The invention now being fully described, it will be apparent to one of ordinary skill in the art that many changes and modifications can be made thereto without departing from the spirit or scope of the appended claims. Modifications that become apparent from the foregoing description and accompanying figures fall within the scope of the claims.